

# Credible Persuasion\*

Xiao Lin<sup>†</sup>

Ce Liu<sup>‡</sup>

June 30, 2023

## Abstract

We propose a new notion of credibility for Bayesian persuasion problems. A disclosure policy is credible if the Sender cannot profit from tampering with her messages while keeping the message distribution unchanged. We show that the credibility of a disclosure policy is equivalent to a cyclical monotonicity condition on its induced distribution over states and actions. We also characterize how credibility restricts the Sender's ability to persuade under different payoff structures. In particular, when the Sender's payoff is state-independent, all disclosure policies are credible. We apply our results to the market for lemons, and show that no useful information can be credibly disclosed by the seller, even though a seller who can commit to her disclosure policy would perfectly reveal her private information to maximize profit.

---

\*We are indebted to Nageeb Ali for his continuing guidance and support. This paper has benefited from the thoughtful and constructive feedback of the editor, Emir Kamenica, and three anonymous referees. We are also grateful for comments and suggestions from Ian Ball, Carl Davidson, Jon Eguia, Henrique de Oliveira, Piotr Dworczak, Alex Frankel, Nima Haghpahan, Rick Harbaugh, Marc Henry, Tetsuya Hoshino, Yuhta Ishii, Navin Kartik, Vijay Krishna, SangMok Lee, George Mailath, Meg Meyer, Moritz Meyer-ter-Vehn, Harry Pei, Daniel Rappoport, Andrew Rhodes, Ron Siegel, Alex Smolin, Juuso Toikka, Rakesh Vohra, Jia Xiang, Takuro Yamashita, and participants at various conferences and seminars. Siqi Li provided excellent research assistance.

<sup>†</sup>Department of Economics, University of Pennsylvania (e-mail: [xiaolin7@upenn.edu](mailto:xiaolin7@upenn.edu)).

<sup>‡</sup>Department of Economics, Michigan State University (e-mail: [celiu@msu.edu](mailto:celiu@msu.edu)).

# Contents

|          |   |           |
|----------|---|-----------|
| <b>1</b> | <b>Introduction</b>                                 | <b>1</b>  |
| <b>2</b> | <b>Model</b>  | <b>7</b>  |
| 2.1      | Setup . . . . .                                     | 7         |
| 2.2      | Stable Outcome Distributions . . . . .              | 9         |
| 2.3      | The Case of State-Independent Preferences . . . . . | 11        |
| 2.4      | When is Credibility Restrictive? . . . . .          | 12        |
| <b>3</b> | <b>Application: The Market for Lemons</b>           | <b>19</b> |
| <b>4</b> | <b>Discussion</b>                                   | <b>22</b> |
| 4.1      | Relationship to Rochet (1987) . . . . .             | 22        |
| 4.2      | Finite-Sample Approximation . . . . .               | 24        |
| 4.3      | Restriction to Pure Strategies . . . . .            | 26        |
| <b>5</b> | <b>Conclusion</b>                                   | <b>28</b> |
|          | <b>References</b>                                   | <b>28</b> |
| <b>A</b> | <b>Appendix</b>                                     | <b>32</b> |
| <b>B</b> | <b>Supplementary Appendix</b>                       | <b>59</b> |

# 1 Introduction

When an informed party (Sender; she) discloses information to persuade her audience (Receiver; he), it is in her interest to convey only messages that steer the outcome in her own favor: schools may want to inflate their grading policies to improve their job placement records; similarly, credit rating agencies may publish higher ratings in exchange for future business. Even when the Sender claims to have adopted a disclosure policy, she may still find it difficult to commit to following its prescriptions, since the adherence to such policies is often impossible to monitor. By contrast, what *is* often publicly observable is the final distribution of the Sender’s messages: students’ grade distributions at many universities are publicly available, and so are the distributions of rating scores from credit rating agencies.

Motivated by this observation, we propose a notion of *credible persuasion*. In contrast to standard Bayesian persuasion, our Sender cannot commit to a disclosure policy; however, to avoid detection, she must keep the final message distribution unchanged when deviating from her disclosure policy. For example, in the context of grade distributions, if a university had announced a disclosure policy that features certain fractions of A’s, B’s, and C’s, it cannot switch to a distribution that assigns every student an A without being detected. Analogously, if a credit rating agency were to tamper with its rating scheme, any resulting change in the overall distribution of ratings may be detected. Our notion of credibility closely adheres to this definition of detectability: we say that a disclosure policy is credible if given how the Receiver reacts to her messages, the Sender has no profitable deviation to any other disclosure policy that has the same message distribution.

Can the Sender persuade the Receiver by using credible disclosure policies? We find that in many settings, no informative disclosure policy is credible. An important case where this effect is exhibited is the market for lemons (Akerlof, 1970). Here, we show that the seller of an asset cannot credibly disclose any useful information to the buyer; this effect arises even though the seller benefits from persuasion when she can fully commit to her disclosure policy. Conversely, we also provide conditions for when the Sender is guaranteed to benefit from credible persuasion so that credibility does not entirely eliminate the scope for persuasion. In general, we show that credibility is characterized by a *cyclical monotonicity* condition, which is analogous to those studied in decision theory and mechanism design (Rochet, 1987).

To illustrate these ideas, consider the following example. A buyer (Receiver) chooses whether to buy a car from a used car seller (Sender). It is common knowledge that 30% of the cars are of *high* quality and the remaining 70% are of *low* quality. For simplicity, suppose that all cars are sold at an exogenously fixed price.<sup>1</sup> The payoffs in this example are in Table 1.

---

<sup>1</sup>In Section 3 we study a competitive market for lemons with endogenous prices, and emerge with similar findings.

|        | Buy | Not Buy |       | Buy | Not Buy |
|--------|-----|---------|-------|-----|---------|
| High   | 2   | 1       | High  | 1   | 0       |
| Low    | 2   | 0       | Low   | -1  | 0       |
| Seller |     |         | Buyer |     |         |

Table 1: Used Car Example Payoffs

The seller always prefers selling a car, but the buyer is willing to purchase if and only if he believes its quality is high with at least 0.5 probability. Conditional on a car being sold, the seller obtains the same payoff regardless of its quality; but when a car is not sold, she receives a higher value from retaining a high-quality car.

As a benchmark, let us first see what the seller achieves if she could commit to a disclosure policy. We depict the optimal disclosure policy in Figure 1. The policy uses two messages, *pass* and *fail*: all high-quality cars pass, along with 3/7 of the low-quality cars; the remaining 4/7 of the low-quality cars receive a failing grade. Conditional on the car passing, the buyer believes that the car is of high quality with probability 0.5, which is just enough to convince him to make the purchase. If a car fails, the buyer believes that the car is of low quality for sure and will refuse to buy. With this disclosure policy, the buyer expects to see the seller pass 60% of the cars and fail the remaining 40%.<sup>2</sup>

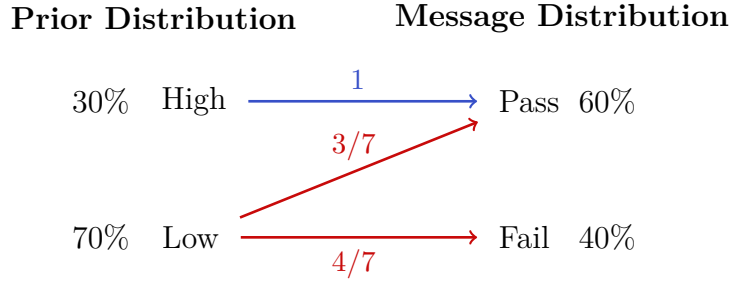


Figure 1: Optimal Commitment Policy

The policy above is optimal for the seller if she can commit to following its prescriptions. But suppose the buyer cannot observe how the seller rates her cars. Instead, the buyer only observes the fraction of cars being passed and failed. In such a setting, the seller can profitably deviate from the above disclosure policy without being detected by the buyer. Specifically, the seller can switch to failing all high-quality cars while adding an equal number of low-quality cars to the passing grade. This disclosure policy, illustrated in Figure 2, induces the same distribution of messages (i.e., 60% pass, 40% fail). Holding fixed the buyer's behavior,

<sup>2</sup>This example, by design, has the same solution as the prosecutor-judge example in Kamenica and Gentzkow (2011).

this deviation is profitable for the seller because she still ends up selling the same number of cars but now is able to retain more high-quality cars. Accordingly, we view the optimal full-commitment policy to be not credible: after having promised to share information according to a disclosure policy, the seller would not find it rational to follow through and would instead profit from an undetectable deviation.

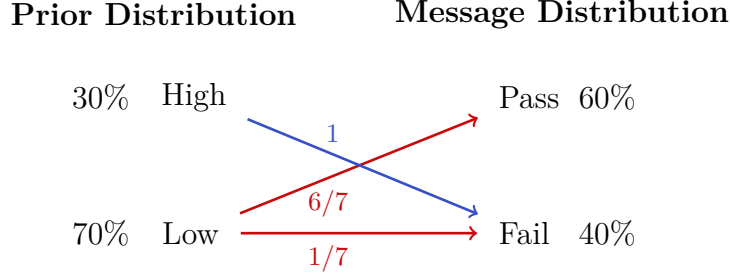


Figure 2: An Undetectable Deviation

More generally, we introduce the following notion of credibility for disclosure policies. Consider a *profile* consisting of the Sender’s disclosure policy and the Receiver’s strategy (mapping messages to actions). We say that a profile is *Receiver incentive compatible* if the Receiver’s strategy best responds to the Sender’s disclosure policy—this requirement is standard in Bayesian persuasion problems. We say that a profile is *credible* if, given the Receiver’s strategy, the Sender has no profitable deviation to any other disclosure policy that induces the same message distribution. Together, credibility and Receiver incentive compatibility require that conditional on the Sender’s message distribution, the Sender and Receiver best respond to each other.<sup>3</sup>

We have just argued that in the used car example, the optimal full-commitment disclosure policy was not credible given the Receiver’s best response. Can any car be sold in a profile that is both credible and Receiver incentive compatible? The answer is no. Note that zero sales is also the outcome when no information is disclosed. In other words, credibility completely shuts down the possibility for useful information transmission.

To see why, suppose towards a contradiction that the buyer purchases a car after observing a message  $m_1$  that is sent with positive probability. By Receiver incentive compatibility, the buyer must believe that the car is of high quality with at least 0.5 probability after observing  $m_1$ . Since  $m_1$  is sent with positive probability, the martingale property of beliefs implies that there must be another message  $m_2$ , also sent with positive probability, that makes the buyer assign less than 0.5 to the car’s quality being high. Necessarily, when the buyer observes the message  $m_2$ , he does not make a purchase. This creates an incentive for the seller to tamper

<sup>3</sup>Our solution-concept is therefore analogous to an equilibrium condition in which the set of feasible deviations for the Sender is to other disclosure policies that induce the same message distribution.

with her disclosure policy: by exchanging some of the good cars being mapped into  $m_1$  with an equal number of bad cars being mapped into  $m_2$ , she can improve her payoff without changing the distribution of messages.

One may wonder if credibility always shuts down communication entirely. The next example features a setting in which the optimal full-commitment disclosure policy is credible. Consider the disclosure problem faced by a school (Sender) and an employer (Receiver).<sup>4</sup> Just as in the used car example, a student’s ability is either *high* with probability 0.3 or *low* with probability 0.7. Payoffs are as shown in Table 2. The employer is willing to hire a student if he believes the student has high ability with at least 0.5 probability. The school would like all its students to be employed, but derives a higher payoff from placing a good student than it does from placing a bad one.

|        | Hire | Not Hire |          | Hire | Not Hire |
|--------|------|----------|----------|------|----------|
| High   | 2    | 0        | High     | 1    | 0        |
| Low    | 1    | 0        | Low      | −1   | 0        |
| School |      |          | Employer |      |          |

Table 2: School Example Payoffs

The school’s optimal full-commitment disclosure policy is identical to the one in the used car example (Figure 1), and so are the employer’s best responses. But unlike the used car example, the school cannot profitably deviate without changing the message distribution.

To see why, note that without changing the message distribution, any deviation must involve passing some low ability students while failing an equal number of high ability students. This would increase the employment of low ability students at the expense of their high ability counterparts, which makes the school worse off. Since the school cannot profit from undetectable deviations, the optimal full-commitment policy is credible. In contrast to the previous example where credibility shuts down all useful communication, the current example shows that credibility sometimes imposes no cost on the Sender relative to persuasion with full commitment.

In the two examples above, credibility has starkly different implications for information transmission. The key difference is that in the used car example, when the car’s quality is higher, the Sender has a weaker incentive to trade while the Receiver’s incentive to trade is stronger; in the school example, by contrast, both the Sender and Receiver have a stronger incentive to trade as the student’s ability increases. Our results formalize this intuition.

<sup>4</sup>See Ostrovsky and Schwarz (2010) for an early study of how schools strategically design their grading policies in a competitive setting.

[Proposition 1](#) shows that when the Sender and Receiver’s preferences have opposite modularities (e.g. when the Sender’s payoff is strictly supermodular and the Receiver’s payoff is submodular), no useful information can be credibly communicated. Even when players’ preferences share the same modularity, the Sender does not always benefit from credible persuasion relative to the no-information benchmark. [Proposition 2](#) and [Proposition 3](#) provide additional conditions that guarantee the Sender does benefit from credible persuasion, as well as conditions under which the optimal full-commitment disclosure policy is credible. [Proposition 4](#) provides a comparative statics result on preference alignment.

Generalizing further, we use optimal transport theory to characterize credibility using a familiar condition from mechanism design and decision theory—cyclical monotonicity. [Theorem 1](#) shows that for every profile of Sender’s disclosure policy and Receiver’s strategy, the credibility of the profile is equivalent to a cyclical monotonicity condition on its induced distribution over states and actions. As is illustrated in the examples above, credibility requires that the Sender cannot benefit from any pairwise swapping in the matching of states and actions. The cyclical monotonicity condition generalizes this idea to cyclical swapping: for every sequence of state-action pairs in the support, the sum of the Sender’s utility should be lower after the matchings of states and actions in this sequence are permuted. In [Section 4.1](#), we discuss the connection of [Theorem 1](#) to [Rochet \(1987\)](#).

Our paper also offers foundations for studying Bayesian persuasion in a number of settings. One example is when the Sender’s payoff is state-independent: in these cases, our results imply that all disclosure policies are credible, so the full-commitment assumption in the Bayesian persuasion approach is nonessential as long as the message distribution is observable. Another example is when the Sender’s payoff is supermodular, in which case all monotone disclosure policies are credible.

The rest of the paper is organized as follows: [Section 2](#) introduces our credibility notion as well as the main results. [Section 3](#) considers an application: in the market for lemons with endogenous prices, we show that the seller cannot credibly disclose any useful information to the buyers, even though full disclosure would maximize the seller’s profit. [Section 4](#) discusses several aspects of our model. [Section 5](#) concludes. All omitted proofs are in [Appendix A](#). The remainder of this introduction places our contribution within the context of the broader literature.

**Related Literature:** Our work contributes to the study of strategic communication. The Bayesian persuasion model in [Kamenica and Gentzkow \(2011\)](#) studies a Sender who can fully

commit to a disclosure policy.<sup>5</sup> By contrast, the cheap-talk approach pioneered by [Crawford and Sobel \(1982\)](#) models a Sender who observes the state privately and, given the Receiver’s strategy, chooses an optimal (sequentially rational) message. The partial-commitment setting that we model is between these two extremes: here, the Sender can commit to a distribution over messages but not the entire disclosure policy.

Our model considers a Sender who can misrepresent her messages as long as the misrepresentation still produces the original message distribution. This contrasts with existing approaches to modeling limited commitment in Bayesian persuasion. One approach, pioneered by [Fréchette, Lizzeri, and Perego \(2021\)](#), [Min \(2021\)](#), and [Lipnowski, Ravid, and Shishkin \(2022\)](#), is to allow the Sender to alter the messages from her chosen disclosure policy with some fixed probability. A different method of modeling limited commitment is to consider settings where the Sender can misreport at a cost.<sup>6</sup> For example, [Guo and Shmaya \(2021\)](#) study a Sender who pays a cost when the posterior beliefs induced by her messages are miscalibrated from their literal meanings; [Nguyen and Tan \(2021\)](#) consider a Sender who can costly revise the messages from her chosen disclosure policy; [Perez-Richet and Skreta \(2021\)](#) consider a Sender who can falsify the state, or input, of the disclosure policy. Another approach, taken in [Libgober \(2022\)](#), is to consider a Sender who publicly chooses some dimension of the signal structure while privately choosing the other dimension. Finally, [Perez-Richet \(2014\)](#), [Hedlund \(2017\)](#), [Koessler and Skreta \(2021\)](#), and [Zapechelnyuk \(2023\)](#) allow the Sender to have private information before choosing the disclosure policy. In these settings, Receiver infers the state through the messages from the disclosure policy as well as the signaling effect of the Sender’s choice of information structures.

The way that we model the Sender’s feasible deviations is closely related to the literature on quota mechanisms, which use message budgets to induce truth-telling; see, for example, [Jackson and Sonnenschein \(2007\)](#), [Matsushima, Miyazaki, and Yagi \(2010\)](#), [Rahman \(2010\)](#), and [Frankel \(2014\)](#). Similar ideas have also been explored in communication games. For example, [Chakraborty and Harbaugh \(2007\)](#) consider multi-issue cheap-talk problems, and study equilibria where the Sender assigns a ranking to each issue. In such equilibria, a message is a complete or partial ordering of all the issues, and any on-path deviation is a different ordering that maintains the same distribution of rankings. [Renault, Solan, and Vieille \(2013\)](#) study repeated cheap-talk models where only messages and the Receiver’s actions are publicly observable. They characterize equilibria in the repeated communication game via a static reporting game where the Sender directly reports her type. The key condition in their characterization requires truthful reporting to be optimal among all reporting strategies that

---

<sup>5</sup>[Brocas and Carrillo \(2007\)](#) and [Rayo and Segal \(2010\)](#) also study optimal disclosure policy in more specific settings.

<sup>6</sup>This approach was initially introduced by [Kartik \(2009\)](#) to study language inflation.



replicate the true type distribution, which is akin to [Rahman \(2010\)](#)’s characterization of implementable direct mechanisms. [Margaria and Smolin \(2018\)](#) use a different approach to study the case where the Sender’s payoff is state-independent, and [Meng \(2021\)](#) provides a unified approach to characterizing the Receiver’s optimal value in these repeated cheap-talk models. [Kuvalekar, Lipnowski, and Ramos \(2021\)](#) study a related model where the Receiver is short-lived, and show that the equilibrium payoffs can be characterized via a static cheap-talk model with capped money burning.

A different strand of the repeated cheap-talk literature studies models where the Receiver can observe feedback about past state realizations. [Best and Quigley \(2020\)](#) consider how coarse feedback of past states can substitute for commitment; [Mathevet, Pearce, and Stacchetti \(2022\)](#) allow for the possibility of non-strategic commitment types; [Pei \(2020\)](#) studies a setting where the Sender has persistent private information about her lying cost.

Finally, our approach to credible persuasion is reminiscent of how [Akbarpour and Li \(2020\)](#) model credible auctions. They study mechanism design problems where the designer’s deviations are “safe” so long as they lead to outcomes that are possible when she is acting honestly, and characterize mechanisms that ensure the designer has no safe and profitable deviations. By contrast, we study persuasion problems where the Sender’s deviations are undetectable if they do not alter the message distribution, and characterize disclosure policies where the Sender has no profitable and undetectable deviations.

## 2 Model

### 2.1 Setup

We consider an environment with a single Sender ( $S$ ; she) and a single Receiver ( $R$ ; he). Both players’ payoffs depend on an unknown state  $\theta \in \Theta$  and the Receiver’s action  $a \in A$ . Both  $\Theta$  and  $A$  are finite sets.<sup>7</sup> The payoff functions are given by  $u_S : \Theta \times A \rightarrow \mathbb{R}$  and  $u_R : \Theta \times A \rightarrow \mathbb{R}$ . Players hold full-support common prior  $\mu_0 \in \Delta(\Theta)$ .

Let  $M$  be a finite message space that contains  $A$ . The Sender chooses an information structure to influence the Receiver’s action. Specifically, an information structure  $\lambda \in \Delta(\Theta \times M)$  is a joint distribution of states and messages, so that the marginal distribution of states agrees with the prior; that is,  $\lambda_\Theta = \mu_0$ .<sup>8</sup> The Receiver chooses an action after observing each message according to a pure strategy  $\sigma : M \rightarrow A$ .<sup>9</sup>

<sup>7</sup>In [Appendix B.1](#), we show that our main characterization result extends to the case where  $\Theta$  and  $A$  are compact Polish spaces.

<sup>8</sup>For a probability measure  $P$  defined on some product space  $X \times Y$ , we use  $P_X$  and  $P_Y$  to denote its marginal distribution on  $X$  and  $Y$ , respectively.

<sup>9</sup>We focus on pure strategies to abstract from the Receiver using randomization to deter the Sender’s

Our interest is in understanding the Sender's incentives to deviate from her information structure, which depends on the Receiver's strategy. To avoid ambiguity, we refer explicitly to pairs of  $(\lambda, \sigma)$ —or *profiles*—that consist of a Sender's information structure and a Receiver's strategy. For each profile  $(\lambda, \sigma)$ , the players' expected payoffs are

$$U_S(\lambda, \sigma) = \sum_{\theta, m} u_S(\theta, \sigma(m)) \lambda(\theta, m) \quad \text{and} \quad U_R(\lambda, \sigma) = \sum_{\theta, m} u_R(\theta, \sigma(m)) \lambda(\theta, m).$$

We consider a setting where the Sender cannot commit to her information structure, and can deviate to another information structure so long as it leaves the final message distribution unchanged. This embodies the notion that the distribution of the Sender's messages is observable, even though it may be difficult to observe exactly how these messages are generated. Formally, if  $\lambda$  is an information structure promised by the Sender, let  $D(\lambda) \equiv \{\lambda' \in \Delta(\Theta \times M) : \lambda'_\Theta = \mu_\Theta, \lambda'_M = \lambda_M\}$  denote the set of information structures that induce the same distribution of messages as  $\lambda$ : these information structures are indistinguishable from  $\lambda$  from the Receiver's perspective. Our credibility notion requires that conditioning on how the Receiver responds to the Sender's messages, no deviation in  $D(\lambda)$  can be profitable for the Sender.

**Definition 1.** A profile  $(\lambda, \sigma)$  is *credible* if

$$\lambda \in \arg \max_{\lambda' \in D(\lambda)} \sum_{\theta, m} u_S(\theta, \sigma(m)) \lambda'(\theta, m). \quad (1)$$

Moreover, the Receiver's strategy is required to be a best response to the Sender's information structure.

**Definition 2.** A profile  $(\lambda, \sigma)$  is *Receiver incentive compatible (R-IC)* if

$$\sigma \in \arg \max_{\sigma' : M \rightarrow A} \sum_{\theta, m} u_R(\theta, \sigma'(m)) \lambda(\theta, m). \quad (2)$$

Together, credibility and R-IC ensure that conditioning on the message distribution of the Sender's information structure, both the Sender and the Receiver best respond to each other. An immediate observation is that there always exists a “babbling” profile  $(\lambda, \sigma)$  that is both credible and R-IC: a degenerated information structure that sends only one message, and a Receiver strategy that best responds to the prior after observing this message.

Note that the credibility notion can be viewed as merely incorporating an additional constraint in the design of information structures. Some of our results focus on Sender optimality, deviations. This restriction is not without loss of generality, though some of our results can be extended to allow Receiver mixing. See [Section 4.3](#) for a more detailed discussion of this assumption.

but the notion can be applied to different design objectives. It is also worth noting that credibility is a constraint that is independent from Receiver incentive compatibility. As a result, our credibility notion can be applied more broadly to settings where the consequences of the Sender’s messages can be specified via an “outcome function.” As an application, we apply our credibility notion to a setting with multiple Receivers in [Section 3](#).

Finally, our credibility notion is motivated by the observability of the Sender’s message distribution, which we model as a restriction on the Sender’s feasible deviations. The observability of message distributions is best understood through a population interpretation of persuasion models,<sup>10</sup> where there is a continuum of objects with types distributed according to  $\mu_0 \in \Delta(\Theta)$ . The Sender’s information structure  $\lambda$  assigns each object a message based on its type, which generates a message distribution  $\lambda_M$ . Working with a continuum population affords us a cleaner exposition by abstracting from sampling variation. In [Section 4.2](#), we consider a finite approximation where the Sender privately observes  $N$  i.i.d. samples from  $\mu_0 \in \Delta(\Theta)$ , and assigns each realization a message  $m \in M$  subject to quotas on message frequencies; the Receiver then chooses an action after observing the Sender’s message. We show that credible and R-IC profiles in our continuum model are approximated by those in the finite-sample model when the sample size  $N$  becomes large.

## 2.2 Stable Outcome Distributions

We characterize credible and Receiver incentive compatible profiles through the induced probability distribution of states and actions. Formally, an *outcome distribution* is a distribution  $\pi \in \Delta(\Theta \times A)$  that satisfies  $\pi_\Theta = \mu_0$ : this is a consistency requirement that stipulates that the marginal distribution of states must conform to the prior. We say an outcome distribution  $\pi$  is induced by a profile  $(\lambda, \sigma)$  if for every  $(\theta, a) \in \Theta \times A$ ,  $\pi(\theta, a) = \lambda(\theta, \sigma^{-1}(a))$ , where  $\sigma^{-1}$  is the inverse mapping of  $\sigma$ . We are interested in characterizing outcome distributions that can be induced by profiles that are both credible and R-IC, and refer to such outcome distributions as stable.

**Definition 3.** *An outcome distribution  $\pi \in \Delta(\Theta \times A)$  is **stable** if it is induced by a profile  $(\lambda, \sigma)$  that is both credible and R-IC.*

Our first result characterizes stable outcome distributions.

**Theorem 1.** *An outcome distribution  $\pi \in \Delta(\Theta \times A)$  is stable if and only if:*

---

<sup>10</sup>For a more detailed discussion of various interpretations of Bayesian persuasion models, see e.g. [Section 2.2 of Kamenica \(2019\)](#).

1.  $\pi$  is  $u_R$ -obedient: for each  $a \in A$  such that  $\pi(\Theta, a) > 0$ ,

$$\sum_{\theta \in \Theta} \pi(\theta, a) u_R(\theta, a) \geq \sum_{\theta \in \Theta} \pi(\theta, a) u_R(\theta, a') \text{ for all } a' \in A.$$

2.  $\pi$  is  $u_S$ -cyclically monotone: for any sequence  $(\theta_1, a_1), \dots, (\theta_n, a_n) \in \text{supp}(\pi)$  where  $a_{n+1} \equiv a_1$ ,

$$\sum_{i=1}^n u_S(\theta_i, a_i) \geq \sum_{i=1}^n u_S(\theta_i, a_{i+1}).$$

The first condition is the standard obedience constraint (Bergemann and Morris, 2016; Taneva, 2019), which specifies that the Receiver finds it incentive compatible to follow the recommended action given the belief that she forms when receiving that recommendation. The second condition, namely  $u_S$ -cyclical monotonicity, is the new constraint that maps directly to our notion of credibility. While both the necessity and sufficiency of  $u_S$ -cyclical monotonicity can be proven by invoking the Kantorovich duality from optimal transport theory, below we outline a direct proof to better illustrate the intuition behind  $u_S$ -cyclical monotonicity. The full version of this proof can be found in Lemma 2 in Appendix A.

Consider an outcome distribution  $\pi$  and a sequence  $(\theta_i, a_i)_{i=1}^n$  in the support of  $\pi$ . For intuition, let us regard  $\pi$  as a direct-recommendation information structure. A “cyclical” deviation in this case consists of subtracting  $\varepsilon$  mass from  $(\theta_i, a_i)$  while adding it to  $(\theta_i, a_{i+1})$  for each  $i = 1, \dots, n$ , where  $a_{n+1} \equiv a_1$ . Each step of this cyclical deviation changes the Sender’s payoff by  $\varepsilon[u_S(\theta_i, a_{i+1}) - u_S(\theta_i, a_i)]$ , so the total change in the Sender’s payoff is

$$\varepsilon \left[ \sum_{i=1}^n u_S(\theta_i, a_{i+1}) - \sum_{i=1}^n u_S(\theta_i, a_i) \right].$$

The cyclical monotonicity condition requires that the Sender can find no profitable cyclical deviations.

To see why this is necessary for credibility, observe that cyclical deviations do not change the distribution of action recommendations, so any such deviation cannot be detected solely on the basis of the distribution of messages. Credibility requires that these undetectable deviations are not profitable, which implies the cyclical monotonicity condition.

For sufficiency, a key observation is that any outcome distribution  $\pi \in \Delta(\Theta \times A)$  can be approximated by a distribution with rational marginals, which can then be normalized and transformed into doubly stochastic matrices. According to the Birkhoff-von Neumann theorem, permutation matrices form the extreme points of all doubly stochastic matrices. In addition, each permutation matrix is equivalent to a cyclical deviation. So in a rough

sense, cyclical deviations are (approximately) the extreme points of all undetectable Sender deviations. It is therefore sufficient to ensure no cyclical deviations are profitable.

The next two results establish additional properties of the credibility notion. [Corollary 1](#) proves the existence of a Sender-optimal credible and R-IC profile, and shows that it needs not involve more than  $\min\{|\Theta|, |A|\}$  messages.

**Corollary 1.** *There exists a Sender-optimal credible and R-IC profile  $(\lambda^*, \sigma^*)$  where  $\lambda^*$  has no more than  $\min\{|\Theta|, |A|\}$  messages.*

The proof idea is analogous to a similar result for optimal persuasion under full commitment. Every outcome distribution can be regarded as a “direct recommendation” information structure that uses less than  $|A|$  messages. Carathéodory’s theorem then delivers the  $|\Theta|$  bound: given any Sender-optimal stable outcome distribution that involves more than  $|\Theta|$  actions (which are treated as messages), applying Carathéodory’s theorem results in an outcome distribution that maintains the same Sender payoff but has a smaller support, which relaxes both the  $u_S$ -cyclical monotonicity condition and the  $u_R$ -obedience condition.

The next corollary ensures that when checking the cyclical monotonicity condition, one can restrict attention to deviations with length  $n \leq \min\{|\Theta|, |A|\}$ . In fact, if there is any profitable cyclical deviation with length longer than  $\min\{|\Theta|, |A|\}$ , we can split it into two shorter cyclical deviations, at least one of which is profitable. As a result, when both  $\Theta$  and  $A$  are finite, [Corollary 2](#) implies that one only needs to check a finite number of deviations.<sup>11</sup>

**Corollary 2.** *An outcome distribution  $\pi$  is  $u_S$ -cyclically monotone if and only if for each sequence  $(\theta_1, a_1), \dots, (\theta_n, a_n) \in \text{supp}(\pi)$  where  $n \leq \min\{|\Theta|, |A|\}$  and  $a_{n+1} \equiv a_1$ , we have*

$$\sum_{i=1}^n u_S(\theta_i, a_i) \geq \sum_{i=1}^n u_S(\theta_i, a_{i+1}).$$

## 2.3 The Case of State-Independent Preferences

If  $u_S(\theta, a)$  is state independent, then  $u_S$ -cyclical monotonicity is automatically satisfied. So we have the following observation.<sup>12</sup>

**Observation 1.** *If  $u_S(\theta, a) = h(a)$  for some  $h : A \rightarrow \mathbb{R}$ , then every outcome distribution that satisfies  $u_R$ -obedience is stable.*

<sup>11</sup>The total number of deviations, even though finite, can in general be quite large, which may make it challenging to verify cyclical monotonicity. In [Section 2.4](#), we impose additional structures on players’ payoffs to gain further tractability.

<sup>12</sup>The same observation holds if  $u(\theta, a) = r(\theta) + h(a)$  for some  $r : \Theta \rightarrow \mathbb{R}$  and  $h : A \rightarrow \mathbb{R}$ , as adding an action-independent nuisance term does not change the Sender’s preferences over outcome distributions given the exogenous prior distribution on  $\Theta$ .

Therefore, in this case, there is no gap between what is achievable by a Sender who can fully commit to an information structure relative to a Sender who can only commit to a distribution of messages.

State-independent payoffs feature in many analyses of communication and persuasion (e.g. Chakraborty and Harbaugh, 2010; Alonso and Câmara, 2016; Lipnowski and Ravid, 2020; Lipnowski, Ravid, and Shishkin, 2021; Gitmez and Molavi, 2022). In these settings, when the Sender is not disclosing information about a population (thus making it difficult to observe the message distribution), generally the optimal full-commitment outcome cannot be achieved. By contrast, our analysis suggests that when one can adopt the population interpretation for Bayesian persuasion models, the Sender can exercise full commitment power by making public the distribution of her messages.

## 2.4 When is Credibility Restrictive?

When the state and action interact in the Sender’s payoff, credibility limits the Sender’s choice of information structures. The goal of this section is to understand how these limits can restrict the Sender’s ability to persuade the Receiver.

In the examples in Section 1, we see that whether the Sender can credibly persuade the Receiver depends crucially on the alignment of their marginal incentives to trade. To understand this logic more generally, we assume that  $\Theta$  and  $A$  are totally ordered sets, which without loss of generality can be assumed to be subsets of  $\mathbb{R}$ . Recall that a payoff function  $u : \Theta \times A \rightarrow \mathbb{R}$  is *supermodular* if for all  $\theta > \theta'$  and  $a > a'$ , we have

$$u(\theta, a) + u(\theta', a') \geq u(\theta, a') + u(\theta', a),$$

and *submodular* if

$$u(\theta, a) + u(\theta', a') \leq u(\theta, a') + u(\theta', a).$$

Furthermore, the function is *strictly supermodular* or *strictly submodular* if the inequalities above are strict for  $\theta > \theta'$  and  $a > a'$ .

The modularity of players’ payoff functions captures how the marginal utility from switching to a higher action varies with the state. This generalizes the marginal incentive to trade in the examples in Section 1: intuitively, the Sender and the Receiver have aligned marginal incentives when both players’ payoff functions share the same modularity, and opposed marginal incentives when their payoff functions have opposite modularities. To fix ideas, we will assume that the Sender’s payoff is supermodular and vary the modularity of the Receiver’s payoff.

We now introduce a lemma that simplifies the  $u_S$ -cyclical monotonicity condition in Theorem 1 when the Sender’s payoff is supermodular. Say that an outcome distribution

$\pi \in \Delta(\Theta \times A)$  is *comonotone* if for all  $(\theta, a), (\theta', a') \in \text{supp}(\pi)$  satisfying  $\theta < \theta'$ , we have  $a \leq a'$ . Comonotonicity requires that the states and the Receiver's actions are positive-assortatively matched in the outcome distribution. The following lemma, whose variant appears in [Rochet \(1987\)](#), shows that  $u_S$ -cyclical monotonicity reduces to comonotonicity when the Sender's preference is supermodular.

**Lemma 1.** *If  $u_S$  is supermodular, then every comonotone outcome distribution is  $u_S$ -cyclically monotone. Furthermore, if  $u_S$  is strictly supermodular, then every  $u_S$ -cyclically monotone outcome distribution is also comonotone.*

Combined with [Theorem 1](#), [Lemma 1](#) implies that when the Sender's preference is strictly supermodular, the credibility of a profile  $(\lambda, \sigma)$  is equivalent to the comonotonicity of its induced outcome distribution. Comonotone outcome distributions have attracted much attention in the persuasion literature in part due to their simplicity and ease of implementation; for example, see [Dworczak and Martini \(2019\)](#), [Goldstein and Leitner \(2018\)](#), [Mensch \(2021\)](#), [Ivanov \(2020\)](#), [Kolotilin \(2018\)](#), and [Kolotilin and Li \(2020\)](#). Our credibility notion provides an additional motivation for focusing on monotone information structures.

**Remark 1.** [Lemma 1](#) is particularly relevant when  $u_S(\theta, a)$  is affine in  $\theta$ : that is, when there exist  $\eta_0(a)$  and  $\eta_1(a)$  such that  $u_S(\theta, a) = \eta_0(a) + \eta_1(a)\theta$  for all  $\theta, a$ . In this case, an outcome distribution  $\pi$  is  $u_S$ -cyclical monotone if and only if for all  $(\theta, a), (\theta', a') \in \text{supp}(\pi)$  with  $\theta < \theta'$ , we have  $\eta_1(a) \leq \eta_1(a')$ . In other words, higher states are matched with actions that lead to higher slope terms in  $u_S(\theta, a)$ . The reason is that we can define an order on  $A$ :  $a' \succeq a$  if and only if  $\eta_1(a') \geq \eta_1(a)$ , so that  $u_S$  is strictly supermodular with respect to such order.<sup>13</sup> The payoff function  $u_S(\theta, a) = \eta_0(a) + \eta_1(a)\theta$  underlies much of the literature on “posterior-mean” problems, which includes several of the papers cited above.

As benchmarks, we will often draw comparisons to what the Sender can achieve when she can fully commit to her information structure, as well as what is achievable when all or no information is disclosed. We say an outcome distribution  $\pi^*$  is an *optimal full-commitment outcome* if it maximizes the Sender's payoff among outcome distributions that satisfy  $u_R$ -obedience. An outcome distribution  $\bar{\pi}$  is a *fully revealing outcome* if the Receiver always chooses a best response to every state; that is,

$$a \in \arg \max_{a' \in A} u_R(\theta, a') \text{ for every } (\theta, a) \in \text{supp}(\bar{\pi}).$$

---

<sup>13</sup>Note that the order  $\succeq$  defined as such may not be antisymmetric. Nevertheless, the proof of [Lemma 1](#) holds as long as  $\succeq$  is complete and transitive. In [Appendix A](#), we prove [Lemma 1](#) without assuming antisymmetry.



Finally, an outcome distribution  $\underline{\pi}$  is a *no-information outcome* if the Receiver always chooses the same action that best responds to the prior belief  $\mu_0$ ; in other words, there exists

$$a^* \in \arg \max_{a \in A} \sum_{\theta \in \Theta} \mu_0(\theta) u_R(\theta, a) \text{ such that } \underline{\pi}_A(a^*) = 1.$$

We say the Sender *benefits from persuasion* if an optimal full-commitment outcome gives the Sender a higher payoff than every no-information outcome. Similarly, we say the Sender *benefits from credible persuasion* if there exists a stable outcome distribution that gives the Sender a strictly higher payoff than every no-information outcome.

**When Credibility Shuts Down Communication:** The next result generalizes the used-car example in [Section 1](#). To simplify the statement of the result, we impose the following regularity assumption on the Receiver's payoff function.

**Assumption 1.** *There exist no distinct  $a, a' \in A$  such that  $u_R(\theta, a) = u_R(\theta, a')$  for all  $\theta \in \Theta$ .*

In other words, from the Receiver's perspective, there are no duplicate actions. This assumption is not without loss, but greatly simplifies the statement of [Proposition 1](#) and [Proposition 2](#) below.

**Proposition 1.** *Under [Assumption 1](#), if  $u_S$  is strictly supermodular and  $u_R$  is submodular, then every stable outcome distribution is a no-information outcome.*

[Proposition 1](#) says that when the players have opposed marginal incentives, credibility completely shuts down information transmission. The logic generalizes what we saw in the used-car example: if two distinct messages resulted in different actions from the Receiver, the Sender and Receiver would have diametrically opposed preferences regarding which action to induce in which state. Therefore whenever R-IC is satisfied, the Sender will have an incentive to deviate to another information structure that swaps states and induces the same marginal distribution of messages.

**When the Sender Benefits from Credible Persuasion:** In light of the school example in [Section 1](#), one might expect credibility to not limit the Sender's ability to persuade when her marginal incentives are aligned with the Receiver's. However, this is false without imposing additional assumptions. For an illustration, consider the following example, in which both the Sender and Receiver have supermodular payoffs. The Sender benefits from persuasion when she can fully commit to her information structure, but no stable outcome distribution can give her a higher payoff than the best no-information outcome.



| $u_S(\theta, a)$ | $a_1$ | $a_2$ | $a_3$ | $a_4$ |
|------------------|-------|-------|-------|-------|
| $\theta = H$     | -1    | 0.75  | 1     | 0     |
| $\theta = L$     | 0     | 0.75  | 0.5   | -1    |

| $u_R(\theta, a)$ | $a_1$ | $a_2$ | $a_3$ | $a_4$ |
|------------------|-------|-------|-------|-------|
| $\theta = H$     | 0     | 0.6   | 0.8   | 1     |
| $\theta = L$     | 1     | 0.8   | 0.6   | 0     |

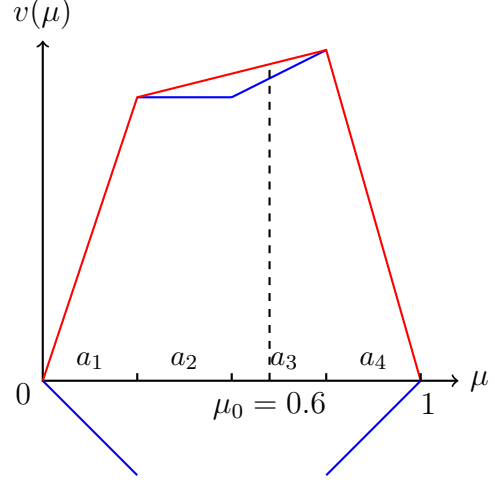


Table 3: Sender and Receiver's payoffs

Figure 3: Concavification

**Example 1.** Suppose  $\Theta = \{H, L\}$  with prior  $\mu_0 = P(\theta = H) = 0.6$  and  $A = \{a_1, a_2, a_3, a_4\}$ . The Sender and Receiver's payoffs are as given in Table 3. Note that both players' payoffs are strictly supermodular. The Receiver's best response is  $a_1$  when  $\mu_0 \in [0, 0.25)$ ,  $a_2$  when  $\mu_0 \in [0.25, 0.5)$ ,  $a_3$  when  $\mu_0 \in [0.5, 0.75)$ , and  $a_4$  when  $\mu_0 \in [0.75, 1]$ ; this leads to the Sender's indirect utility function (blue) and its concave envelope (red) depicted in Figure 3. From Kamenica and Gentzkow (2011), the red line represents the Sender's optimal value under full commitment. It is clear that at  $\mu_0 = 0.6$ , the Sender strictly benefits from persuasion if she can fully commit to her information structure.

However, no stable outcome distribution can make the Sender better off than the no-information outcome. To see why, first note that according to Corollary 1, it is without loss to look for Sender-optimal credible and R-IC profiles that induce only two posterior beliefs  $\mu_1 < \mu_2$ . Now consider the Receiver's actions induced by these two posteriors. By Lemma 1, at most one of these actions can be matched with more than one state, for otherwise the outcome distribution would not be comonotone. So at most one of the actions can be induced by interior posterior beliefs. However, it is clear from Figure 3 that in order for the Sender to benefit from using only two posteriors, she must induce both  $a_2$  and  $a_3$ , both of which can only happen when the Receiver holds interior beliefs. As a result, no credible and R-IC profiles can make the Sender better off.

Example 1 above shows that besides the co-modularity of preferences, additional conditions are needed in order to ensure the Sender can benefit from credible persuasion. Proposition 2 below offers several such conditions.

Let  $A^\circ \equiv \{a \in A : a \in \arg \max_{a'} \sum_{\theta} \mu(\theta) u_R(\theta, a') \text{ for some } \mu \in \Delta(\Theta)\}$  denote the set of

actions that are best responses to some belief of the Receiver; clearly, actions that are not in  $A^\circ$  would never be played by the Receiver in any R-IC profile, and can without loss be discarded from the action set  $A$ . Let  $\bar{a} \equiv \max A^\circ$  and  $\underline{a} \equiv \min A^\circ$  denote the highest and lowest actions in  $A^\circ$ , and let  $\bar{\theta} \equiv \max \Theta$  and  $\underline{\theta} \equiv \min \Theta$  denote the highest and lowest states.

**Proposition 2.** *Suppose both  $u_S$  and  $u_R$  are supermodular, and [Assumption 1](#) holds, then*

1. *If the highest action is dominant for the Sender, that is, if  $u_S(\theta, \bar{a}) > u_S(\theta, a)$  for all  $\theta$  and  $a \in A^\circ \setminus \{\bar{a}\}$ , then for generic priors,<sup>14</sup> the Sender benefits from credible persuasion as long as she benefits from persuasion.*
2. *If the Sender favors extreme actions in extreme states, that is, if  $u_S(\bar{\theta}, \bar{a}) > u_S(\bar{\theta}, a)$  for all  $a \neq \bar{a}$  and  $u_S(\underline{\theta}, \underline{a}) > u_S(\underline{\theta}, a)$  for all  $a \neq \underline{a}$ , then for generic priors, the Sender benefits from credible persuasion.*
3. *If the Sender is strictly better off from a fully revealing outcome than from every no-information outcome, then the Sender benefits from credible persuasion.*

The first condition in [Proposition 2](#) is satisfied in settings like the school example, where the school and the employer's preferences are both supermodular, and the school would always want to place a student regardless of the student's ability. The second condition is applicable in environments where both parties have agreement on extreme states. For example, both doctors and patients favor aggressive treatment if the patient's condition is severe, and both favor no treatment if the patient is healthy, but they might disagree in intermediate cases. Lastly, a special case of the third condition is quadratic-loss preferences as commonly used in models of strategic communication (e.g. [Crawford and Sobel, 1982](#)).<sup>15</sup> However, note that the conditions in [Proposition 2](#) do not guarantee the Sender her optimal full-commitment payoff. In [Appendix B.3](#), we provide an example satisfying the first condition in [Proposition 2](#). The Sender in this example can benefit from credible persuasion, but is unable to achieve the optimal full-commitment payoff.

The first two parts of [Proposition 2](#) are based on belief splitting. Let us briefly describe the proof for the first condition; the proof for the second part follows similar arguments. Note that if  $\bar{a}$  is a dominant action for the Sender, and the Sender can benefit from persuasion (under full commitment), then  $\bar{a}$  must not already be a best response for the Receiver under the prior  $\mu_0$ . The Sender can then split the prior into a point mass posterior  $\delta_{\bar{\theta}}$  and some

<sup>14</sup>Formally, by generic we mean that the result holds under a set of priors  $T \subset \Delta(\Theta)$  that is open, dense, and has full Lebesgue measure.

<sup>15</sup>The model in this section has finite action spaces, so we need to additionally assume that the action space is rich enough such that the Sender's indirect utility function approximates the one under a continuous action space.

other posterior  $\tilde{\mu}$  that is close to  $\mu_0$ . At  $\delta_{\bar{\theta}}$ , the Receiver is induced to choose  $\bar{a}$  since his payoff is supermodular. In addition, for generic priors the Receiver's best response to  $\tilde{\mu}$  remains the same as his best response to  $\mu_0$ . The Sender benefits from this belief-splitting since the same action is still played most of the time, but in addition her favorite action is now played with positive probability. Moreover, the resulting outcome distribution matches higher states with higher actions, so it is stable due to the supermodularity of  $u_S$  and [Lemma 1](#).

The third part of [Proposition 2](#) follows because the fully revealing outcome distribution is always credible when both players' preferences are supermodular. The intuition of this result is most transparent when the Sender's payoff is strictly supermodular. Consider  $(\theta, a)$  and  $(\theta', a')$  in the support of a fully revealing outcome distribution  $\pi$ , so  $a$  and  $a'$  best respond to  $\theta$  and  $\theta'$ , respectively. From [Topkis \(2011\)](#), it follows that  $a \geq a'$  if  $\theta > \theta'$ . Therefore,  $\pi$  is comonotone and satisfies  $u_S$ -cyclical monotonicity by [Lemma 1](#). By construction,  $\pi$  also satisfies  $u_R$ -obedience, so  $\pi$  is stable by [Theorem 1](#). This result is closely related to [Theorem 1](#) and [Theorem 2](#) of [Chakraborty and Harbaugh \(2007\)](#). They show that in multi-issue cheap-talk problems, truthfully revealing the rankings of the issues is an equilibrium under supermodular preferences; in addition, when the number of issues grows to infinity, revealing their rankings is asymptotically equivalent to revealing their values. The credibility of the fully revealing outcome can therefore be viewed as the limit of a rank revealing equilibrium in [Chakraborty and Harbaugh \(2007\)](#).

**When Credibility Imposes No Cost to the Sender:** In [Observation 1](#), we see that when the Sender's payoff is additively separable, credibility does not restrict the set of stable outcomes. [Proposition 3](#) below provides a condition which guarantees that credibility imposes no loss on the Sender's optimal value, even when credibility does restrict the set of stable outcomes.

**Proposition 3.** *Suppose  $|A| = 2$ . If both  $u_S$  and  $u_R$  are supermodular, then at least one optimal full-commitment outcome is stable; if in addition  $u_S$  is strictly supermodular, then every optimal full-commitment outcome is stable.*

[Proposition 3](#) says that in settings where both players have supermodular payoffs and the Receiver faces a binary decision, such as "accept" or "reject", then credibility imposes no cost to the Sender. This result follows from combining our [Theorem 1](#) and [Lemma 1](#) with [Theorem 1](#) in [Mensch \(2021\)](#). He shows that under the assumptions in our [Proposition 3](#), there exists an optimal full-commitment outcome that is comonotone. The intuition is that for any outcome distribution  $\pi$  that is  $u_R$ -obedient but not comonotone, the Sender can weakly improve her payoff by swapping the non-comonotone pairs in the support of  $\pi$ , so that they become matched assortatively. Such swapping also benefits the Receiver due to

the supermodularity of  $u_R$ , so  $u_R$ -obedience remains satisfied. As a result, the Sender can always transform a non-comonotone outcome distribution into one that is comonotone without violating  $u_R$ -obedience, while weakly improving her own payoff. Therefore, there must be an optimal full-commitment outcome that is comonotone, which is also stable by [Theorem 1](#) and [Lemma 1](#).

**Comparative Statics:** Our analysis thus far demonstrates that the mode of preference alignment plays a crucial role in determining the scope of credible persuasion. In this section, we provide a comparative statics result relating the Sender's optimal credible-persuasion payoff to the degree of preference alignment.

In order to measure the Sender's utility on a constant scale, we will keep the Sender's payoff function unchanged and adjust only the Receiver's payoff.<sup>16</sup> Following Section IV of [Kamenica and Gentzkow \(2011\)](#), we say preferences  $(u_S, u'_R)$  are more aligned than  $(u_S, u_R)$  if for any  $a \in A$  and any  $\mu \in \Delta(\Theta)$ ,

$$E_\mu[u_S(\theta, \hat{a}(\mu))] \geq E_\mu[u_S(\theta, a)] \Rightarrow E_\mu[u_S(\theta, \hat{a}'(\mu))] \geq E_\mu[u_S(\theta, a)],$$

where  $\hat{a}(\mu) \in \arg \max_{a \in A} \sum_\theta \mu(\theta) u_R(\theta, a)$  and  $\hat{a}'(\mu) \in \arg \max_{a \in A} \sum_\theta \mu(\theta) u'_R(\theta, a)$  denote the Receiver's best response function, with ties broken in the Sender's favor.

The following result shows that when payoffs are supermodular and preferences become more aligned, the Sender is guaranteed a higher payoff from credible persuasion.

**Proposition 4.** *Suppose  $u_S$ ,  $u_R$ , and  $u'_R$  are strictly supermodular payoff functions. If in addition the preferences  $(u_S, u'_R)$  are more aligned than  $(u_S, u_R)$ , then under  $(u_S, u'_R)$  the Sender obtains a higher payoff from the Sender-optimal stable outcome distribution compared to under  $(u_S, u_R)$ .*

To prove [Proposition 4](#), we take an optimal stable outcome distribution  $\pi$  under the less aligned preferences  $(u_S, u_R)$ , and show that when this same  $\pi$  is used as an information structure under the more aligned preferences  $(u_S, u'_R)$ , it induces a stable outcome distribution that offers the Sender a superior payoff. Specifically, consider the outcome distribution  $\pi'$  induced by the Receiver choosing the Sender-favored best responses to  $\pi$  under  $(u_S, u'_R)$ . Since  $\pi$  is a stable outcome distribution under  $(u_S, u_R)$  and  $u_S$  is supermodular, it follows that  $\pi$  must be comonotone; this combined with the fact that  $u'_R$  is supermodular implies

---

<sup>16</sup>In fact, the change of  $u_S$  would only affect the Sender's optimal credible-persuasion payoff through a scaling effect: according to [Theorem 1](#) and [Lemma 1](#), credibility is equivalent to the outcome distribution being comonotone as long as  $u_S$  is strictly supermodular. So under our maintained assumptions on payoff functions, the set of stable outcome distributions would be unaffected by modifications in the Sender's payoff function.

that  $\pi'$  is also comonotone, and therefore stable under  $(u_S, u'_R)$ . Moreover, as  $(u_S, u'_R)$  is more aligned than  $(u_S, u_R)$ , following each message from  $\pi$ , the Receiver's chosen action in  $\pi'$  is more favorable to the Sender than the recommended action from  $\pi$ . The Sender obtains a higher payoff from  $\pi'$  compared to  $\pi$ , and therefore must be better off under  $(u_S, u'_R)$  than  $(u_S, u_R)$ .

Finally, we provide an example of a class of preferences that meets the requirements of [Proposition 4](#).

**Example 2.** Let  $u_S(\theta, a)$  be a strictly supermodular payoff function; in addition, assume that  $u_S$  favors higher actions:  $u_S(\theta, a') \geq u_S(\theta, a)$  for all  $\theta$  and  $a' \geq a$ . Let  $\{u_R^\kappa\}_{\kappa \in K}$  denote a collection of Receiver's payoff functions defined by  $u_R^\kappa(\theta, a) \equiv w(\theta, a, \kappa)$ , where  $w : \Theta \times A \times K \rightarrow \mathbb{R}$  is a strictly supermodular function, and  $K \subseteq \mathbb{R}$  represents a parameter space.

It's straightforward to see that for each  $\kappa \in K$ , the Receiver payoff function  $u_R^\kappa : \Theta \times A \rightarrow \mathbb{R}$  is strictly supermodular. Furthermore, preferences  $(u_S, u_R^{\kappa'})$  are more aligned than  $(u_S, u_R^\kappa)$  whenever  $\kappa' \geq \kappa$ . To see why, for each  $\kappa \in K$  and  $\mu \in \Delta(\Theta)$ , let  $\hat{a}^\kappa(\mu) \equiv \max\{\arg \max_{a \in A} \sum_{\theta} \mu(\theta) u_R^\kappa(\theta, a)\}$  denote the Receiver's highest best response to  $\mu$  when the payoff function is  $u_R^\kappa$  (note that since the Sender favors higher actions, selecting the highest best response is equivalent to breaking ties in the Sender's favor). By Lemma 2.8.1 of [Topkis \(2011\)](#),  $\hat{a}^{\kappa'}(\mu) \geq \hat{a}^\kappa(\mu)$  for  $\kappa' \geq \kappa$ . Since the Sender favors higher actions, for any  $a \in A$ ,  $\mu \in \Delta(\Theta)$ , and  $\kappa' \geq \kappa$ , we have

$$E_\mu[u_S(\theta, \hat{a}^\kappa(\mu))] \geq E_\mu[u_S(\theta, a)] \Rightarrow E_\mu[u_S(\theta, \hat{a}^{\kappa'}(\mu))] \geq E_\mu[u_S(\theta, a)].$$

This implies that  $(u_S, u_R^{\kappa'})$  are more aligned than  $(u_S, u_R^\kappa)$  whenever  $\kappa' \geq \kappa$ . So according to [Proposition 4](#), the Sender obtains a higher payoff from the Sender-optimal stable outcome distribution under  $(u_S, u_R^{\kappa'})$  than from that under  $(u_S, u_R^\kappa)$ .<sup>17</sup>

### 3 Application: The Market for Lemons

A classic insight from [Akerlof \(1970\)](#) is that in markets with asymmetric information, adverse selection can lead to substantial efficiency loss. In practice, buyers and sellers often rely on warranty or third-party certification to overcome this inefficiency. A seemingly more direct solution to their predicament is for the seller to fully reveal her private information, so that there is no information asymmetry between players. In this section, however, we show that

<sup>17</sup>Note also that the following variant of the alignment notion in [Gentzkow and Kamenica \(2017\)](#) is a further special case of this class of preferences:  $u_S(\theta, a) = f(\theta, a)$ ,  $u_R^\kappa(\theta, a) = f(\theta, a) + \kappa g(\theta, a)$  with  $\kappa \in [0, \infty)$ , where both  $f$  and  $g$  are strictly supermodular and  $f(\theta, a') \geq f(\theta, a)$  for all  $\theta$  and  $a' \geq a$ .

this apparently easy fix to the adverse selection problem relies on unrealistic assumptions on the seller's ability to commit. Indeed, we show that any information disclosure that improves efficiency cannot be credible.

To fix ideas, we adapt the formulation in [Mas-Colell, Whinston, and Green \(1995\)](#) and consider a seller who values an asset she owns (say, a car) at  $\theta \in \Theta \subseteq [0, 1]$ ; two buyers (1 and 2) both value the car at  $v(\theta)$  which is weakly increasing in  $\theta$ . Buyers share a common prior belief  $\mu_0 \in \Delta(\Theta)$ . We assume  $v(\theta) > \theta$  for all  $\theta \in \Theta$  so there is common knowledge of gain from trade. Moreover, we assume  $E_{\mu_0}[v(\theta)] < 1$  so that without information disclosure, some cars will not be traded due to adverse selection. Below we first describe the base game without information disclosure, then augment the base game to allow the seller to choose an information structure to influence the buyers' beliefs.

**The Base Game  $G$ :** The seller and the buyers move simultaneously. The seller learns her value and chooses an ask price  $a_s \in A_s = [0, v(1)]$ ; each buyer  $i = 1, 2$  chooses a bid  $b_i \in A_i = [0, v(1)]$ . If the ask price is lower than or equal to the highest bid, the car is sold at the highest bid to the winning buyer, and ties are broken evenly. If the ask price is higher than the highest bid, the seller keeps the car and receives the reserve value  $\theta$ , while both buyers get 0. More formally, the seller's payoff function is

$$u_s(\theta, a_s, b_1, b_2) = \begin{cases} \max\{b_1, b_2\} & \text{if } a_s \leq \max\{b_1, b_2\} \\ \theta & \text{if } a_s > \max\{b_1, b_2\} \end{cases}$$

and buyer  $i$ 's payoff is

$$u_i(\theta, a_s, b_1, b_2) = \begin{cases} v(\theta) - b_i & \text{if } b_i > b_{-i} \text{ and } b_i \geq a_s \\ \frac{1}{2}[v(\theta) - b_i] & \text{if } b_i = b_{-i} \text{ and } b_i \geq a_s \\ 0 & \text{otherwise.} \end{cases}$$

**The Game with Disclosure:** Let  $M$  be the set of messages, which we assume is a Polish space. Before the base game is played, the seller chooses an information structure  $\lambda$  to publicly disclose information to the buyers.<sup>18</sup> Together the information structure  $\lambda$  and the base game  $G$  define a Bayesian game  $\langle G, \lambda \rangle$ . Every message  $m$  from the information structure  $\lambda$  induces a posterior belief  $\mu_m \equiv \lambda(\cdot|m) \in \Delta(\Theta)$  for the buyers. The buyers  $i = 1, 2$  choose their respective bids  $\beta_i(m)$ , while the seller chooses an ask price  $\alpha_s(\theta, m)$ . We restrict attention to

---

<sup>18</sup>In our setting,  $\lambda$  determines only the buyers' information structure, and the seller is perfectly informed about  $\theta$ . That is, the seller cannot prevent herself from learning the true quality of the car. This differs from [Kartik and Zhong \(2019\)](#), who fully characterize payoffs in the market for lemons under all possible information structures.

Bayesian Nash equilibria where the seller plays her *weakly dominant* strategy  $\alpha_S(\theta, m) = \theta$ , and buyers play pure strategies. As we show in [Lemma 6](#), such equilibria exist in  $\langle G, \lambda \rangle$  for every  $\lambda$ . These equilibria also give rise to the familiar fixed-point characterization of equilibrium price: buyers' bids satisfy

$$\max \{ \beta_1(m), \beta_2(m) \} = E_{\mu_m} [v(\theta) | \theta \leq \max \{ \beta_1(m), \beta_2(m) \}].$$

The trading game above differs from the Sender-Receiver setting in [Section 2](#) in two ways: first, the Sender in the current setting publicly discloses information to multiple Receivers; second, in addition to the Receivers, the Sender also chooses an action (ask price) after observing the realization of the information structure. Nevertheless, the notion of stable outcome distribution extends to the current setting. In particular, the credibility notion is based on the same idea that the Sender cannot profitably deviate to a different information structure without changing the message distribution. The Receiver incentive compatible (R-IC) condition, meanwhile, is replaced by a new IC condition that asks both the Sender and Receivers to play according to a Bayesian Nash equilibrium in  $\langle G, \lambda \rangle$ . As mentioned above, in the market for lemons we will focus on a special class of Bayesian Nash equilibria in the game  $\langle G, \lambda \rangle$  where the seller plays her *weakly dominant strategy*  $\alpha_S(\theta, m) = \theta$ , and the buyers do not mix. We will call such profiles  $(\lambda, \sigma)$  WD-IC to distinguish from the weaker IC requirement. The formal discussion of our credibility notion in this multiple-Receiver setting is notationally cumbersome, and is deferred to [Appendix B.2](#).

Next we state our result, discuss its implications, and provide intuition for its proof. As a benchmark, fix an arbitrary message  $m_0 \in M$ , and let  $\lambda_0 \equiv \mu_0 \times \delta_{m_0}$  be a null information structure. Let  $R_0$  denote the supremum of the seller's payoffs among profiles  $(\lambda_0, \sigma)$  that are WD-IC, so  $R_0$  represents the highest equilibrium payoff the seller can achieve when providing no information.

**Proposition 5.** *Under every credible and WD-IC profile, the seller's payoff is no more than  $R_0$ .*

[Proposition 5](#) implies that any information that can be credibly disclosed is not going to improve the seller's payoff compared to the no-information benchmark. This is in sharp contrast to the full-commitment case, where the seller would like to fully reveal the car's quality, and all car types  $\theta$  are sold at  $v(\theta)$ , which would allow the seller to capture all surplus from trade.

Let us describe the intuition behind the proof for [Proposition 5](#).<sup>19</sup> For each message  $m$  from the seller's information structure  $\lambda$ , let  $\Theta(m)$  denote the support of the buyer's posterior

---

<sup>19</sup>While the message of [Proposition 5](#) is reminiscent of [Proposition 1](#), it requires a different proof since the seller has a private action, so [Theorem 1](#) does not apply. Instead of working with the outcome distribution



belief after observing  $m$ . A key step in proving [Proposition 5](#) is to show that there exists a common trading threshold  $\tau$  such that for each message  $m$ , a car of quality  $\theta \in \Theta(m)$  is traded if and only if  $\theta \leq \tau$ . To see why, suppose towards a contradiction that the trading threshold in message  $m$  is higher than the threshold in another message  $m'$ . We show in the proof that the seller would then have a profitable deviation by swapping some of the cars slightly below the higher threshold in message  $m$  with an equal amount of cars from  $m'$  that are of worse quality.<sup>20</sup> Because this deviation does not change the seller's message distribution, it is also undetectable. Therefore, credibility demands a common threshold  $\tau$  that applies across messages. Given this common threshold  $\tau$ , we then apply Tarski's fixed-point theorem to show that when no information is disclosed, there is an equilibrium that features a higher trading threshold  $\tau' \geq \tau$ . Since a higher threshold means more cars are being traded, which in turn increases the seller's payoff, the seller's payoff under every stable outcome is therefore weakly worse than her payoff from a no-information outcome, and this proves our result.

## 4 Discussion

### 4.1 Relationship to [Rochet \(1987\)](#)

The  $u_S$ -cyclical monotonicity condition in our characterization closely resembles the cyclical monotonicity condition for implementing transfers in [Rochet \(1987\)](#). The reader might wonder why cyclical monotonicity arises in our setting despite the lack of transfers. The connection is best summarized by the following three equivalent conditions from optimal transport theory (see, for example, Theorem 5.10 of [Villani \(2008\)](#)).

**Kantorovich Duality.** *Suppose  $X$  and  $Y$  are both finite sets, and  $u : X \times Y \rightarrow \mathbb{R}$  is a real-valued function. Let  $\mu$  be a probability measure on  $X$  and  $\nu$  be a probability measure on  $Y$ , and  $\Pi(\mu, \nu)$  be the set of probability measures on  $X \times Y$  such that the marginals on  $X$  and  $Y$  are  $\mu$  and  $\nu$ , respectively. Then for any  $\pi^* \in \Pi(\mu, \nu)$ , the following three statements are equivalent:*

1.  $\pi^* \in \arg \max_{\pi \in \Pi(\mu, \nu)} \sum_{x,y} \pi(x, y) u(x, y);$

---

$\pi \in \Delta(\Theta \times A)$ , here we apply the cyclical monotonicity characterization directly to the seller's information structure  $\lambda \in \Delta(\Theta \times M)$  by invoking an optimal transport result from [Beiglöck, Goldstern, Maresch, and Schachermayer \(2009\)](#).

<sup>20</sup>This deviation is profitable because it allows the seller to replace the higher-quality cars traded in  $m$  with the lower-quality, untraded cars in  $m'$ . After this swapping, the lower-quality cars are now sold at the price for the higher-quality cars in  $m$ , while the higher-quality cars are now retained by the seller in  $m'$ .



2.  $\pi^*$  is  $u$ -cyclically monotone. That is, for any  $n$  and  $(x_1, y_1), \dots, (x_n, y_n) \in \text{supp}(\pi^*)$ ,

$$\sum_{i=1}^n u(x_i, y_i) \geq \sum_{i=1}^n u(x_i, y_{i+1}).$$

3. There exists  $\psi : Y \rightarrow \mathbb{R}$  such that for any  $(x, y) \in \text{supp}(\pi^*)$  and any  $y' \in Y$ ,<sup>21</sup>

$$u(x, y) - \psi(y) \geq u(x, y') - \psi(y').$$

Our [Theorem 1](#) builds on the equivalence between 1 and 2 in the Kantorovich duality theorem above to show the equivalence between credibility and  $u_S$ -cyclical monotonicity.

[Rochet \(1987\)](#)'s classic result on implementation with transfers follows from the equivalence between 2 and 3. To see this, consider a principal-agent problem where the agent's private type space is  $\Theta$  with full-support prior  $\mu_0$ , and the principal's action space is  $A$ . The agent's payoff is  $u(\theta, a) - t$ , where  $t$  is the transfer she makes to the principal. Given an allocation rule  $q : \Theta \rightarrow A$ , let  $v_q(\theta, \theta') \equiv u(\theta, q(\theta'))$  denote the payoff that a type- $\theta$  agent obtains from the allocation intended for type  $\theta'$ . Let  $X = Y = \Theta$  and  $\mu = \nu = \mu_0$  in the Kantorovich duality theorem above, and consider the distribution  $\pi^* \in \Pi(\mu, \nu)$  defined by

$$\pi^*(\theta, \theta') = \begin{cases} \mu_0(\theta) & \text{if } \theta = \theta' \\ 0 & \text{otherwise} \end{cases}$$

By the equivalence of 2 and 3 in the Kantorovich duality theorem,  $\pi^*$  is  $v_q$ -cyclically monotone if and only if there exists  $\psi : \Theta \rightarrow \mathbb{R}$  such that for all  $\theta, \theta' \in \Theta$ ,  $v_q(\theta, \theta) - \psi(\theta) \geq v_q(\theta, \theta') - \psi(\theta')$ . That is,

$$u(\theta, q(\theta)) - \psi(\theta) \geq u(\theta, q(\theta')) - \psi(\theta'),$$

so the allocation rule  $q$  can be implemented by the transfer rule  $\psi : \Theta \rightarrow \mathbb{R}$ . The  $v_q$ -cyclical monotonicity condition says that for every sequence  $\theta_1, \dots, \theta_n \in \Theta$  with  $\theta_{n+1} \equiv \theta_1$ ,

$$\sum_{i=1}^n u(\theta_i, q(\theta_i)) \geq \sum_{i=1}^n u(\theta_i, q(\theta_{i+1})).$$

This is exactly the cyclical monotonicity condition in [Rochet \(1987\)](#).

When  $X = \Theta$  is interpreted as the set of an agent's true types and  $Y = \Theta$  interpreted as the set of reported types, the distribution  $\pi^*$  constructed in the previous paragraph can be

---

<sup>21</sup>This statement can also be equivalently written as: there exists  $\phi : X \rightarrow \mathbb{R}$  and  $\psi : Y \rightarrow \mathbb{R}$ , such that  $\phi(x) + \psi(y) \geq u(x, y)$  for all  $x$  and  $y$ , with equality for  $(x, y)$  in the support of  $\pi^*$ .

interpreted as the agent's truthful reporting strategy. Based on this interpretation, [Rahman \(2010\)](#) uses the duality between 1 and 3 to show that the incentive compatibility of truthful reporting subject to quota constraints is equivalent to implementability with transfers.

## 4.2 Finite-Sample Approximation

As discussed in [Section 2.1](#), we interpret our model as one where the Sender designs an information structure that assigns scores to a large population of realized  $\theta$ 's; in particular, our model abstracts away from sampling variation, so there is no uncertainty in the population's realized type distribution. In this section we explicitly allow sampling variation by considering a finite-sample model where the Sender observes  $N$  random i.i.d. draws from  $\Theta$ , and assigns each realized  $\theta$  a message  $m \in M$ , subject to certain message quotas—in particular, these message quotas substitute for the Sender's commitment to message distributions in the continuum model. We will show that credible and R-IC profiles in our continuum model are approximated by credible and R-IC profiles in the discrete model when the sample size is large.

Consider a finite i.i.d sample of size  $N$  drawn from the type space  $\Theta$  according to the prior distribution  $\mu_0$ . The set of all possible empirical distributions over  $\Theta$  in this  $N$ -sample captures the sampling variations in the realized type distribution, and can be written as

$$\mathcal{F}_\Theta^N = \left\{ f/N : f \in \mathbb{N}^{|\Theta|}, \sum_{\theta \in \Theta} f(\theta) = N \right\}.$$

The Sender assigns each realized  $\theta$  in the  $N$ -sample a message  $m \in M$ , which leads to an  $N$ -sample of messages. Let

$$\mathcal{F}_M^N = \left\{ f/N : f \in \mathbb{N}^{|M|}, \sum_{m \in M} f(m) = N \right\}$$

denote the set of  $N$ -sample empirical distributions over messages. Lastly, for a pair of state and message distributions  $(\mu, \nu)$ , let

$$X^N(\mu, \nu) = \left\{ f/N : f \in \mathbb{N}^{|\Theta| \times |M|}, \sum_{\theta} f(\theta, \cdot) = N\nu(\cdot), \sum_m f(\cdot, m) = N\mu(\cdot) \right\}$$

denote the set of  $N$ -sample empirical joint distributions over states and messages that have marginals  $\mu$  and  $\nu$ . Notice that  $X^N(\mu, \nu) \neq \emptyset$  if and only if  $\mu \in \mathcal{F}_\Theta^N$  and  $\nu \in \mathcal{F}_M^N$ .<sup>22</sup>

---

<sup>22</sup>Notice that for any  $f \in \mathbb{N}^{|\Theta| \times |M|}$ , the sum of any row or column has to be integer, so  $X^N(\mu, \nu) = \emptyset$  if either  $\mu \notin \mathcal{F}_\Theta^N$  or  $\nu \notin \mathcal{F}_M^N$ . On the other hand if  $\mu \in \mathcal{F}_\Theta^N$  and  $\nu \in \mathcal{F}_M^N$ , a  $\lambda \in X^N(\mu, \nu)$  can be constructed from the so-called Northwest corner rule.

Let us now define the  $N$ -sample analogue of credible and R-IC profiles. We consider a Sender who assigns a message  $m \in M$  to each realized  $\theta \in \Theta$  subject to a message quota  $\nu^N \in \mathcal{F}_M^N$ . An  $N$ -sample profile is therefore a triple  $(\nu^N, \phi^N, \sigma^N)$ , where  $\phi^N : \mathcal{F}_\Theta^N \rightarrow \Delta(\Theta \times M)$  is a Sender's strategy that takes every realized empirical distribution over states  $\mu^N \in \mathcal{F}_\Theta^N$  to a joint distribution  $\phi^N(\mu^N) \in X^N(\mu^N, \nu^N)$ ; meanwhile,  $\sigma^N : M \rightarrow A$  is a Receiver's strategy that assigns an action to each observed message.<sup>23</sup>

The definitions of Sender credibility and Receiver incentive compatibility in the  $N$ -sample setting mirror those in our continuum model. In particular, we say an  $N$ -sample profile  $(\nu^N, \phi^N, \sigma^N)$  is credible if for each realized empirical distribution over  $\Theta$ , the Sender always chooses an optimal assignment of messages subject to the message quotas specified in  $\nu^N$ :  $(\nu^N, \phi^N, \sigma^N)$  is credible if for every  $\mu^N \in \mathcal{F}_\Theta^N$ ,

$$\phi(\mu^N) \in \arg \max_{\lambda^N \in X(\mu^N, \nu^N)} \sum_{\theta, m} \lambda^N(\theta, m) u_S(\theta, \sigma^N(m)).$$

We say the  $N$ -sample profile  $(\nu^N, \phi^N, \sigma^N)$  is Receiver incentive compatible (R-IC) if  $\sigma^N$  best-responds to the Sender's strategy  $\phi^N$ . In particular, let  $P^N$  denote the probability distribution over  $\mathcal{F}_\Theta^N$  induced by i.i.d. draws from the prior distribution  $\mu_0 \in \Delta(\Theta)$ , and let  $\phi^N(\theta, m | \mu^N)$  be the probability assigned to  $(\theta, m)$  in the joint distribution  $\phi^N(\mu^N)$  chosen by the Sender. The profile  $(\nu^N, \phi^N, \sigma^N)$  is R-IC if

$$\sigma^N \in \arg \max_{\sigma' : M \rightarrow A} \sum_{\mu^N \in \mathcal{F}_\Theta^N} P^N(\mu^N) \sum_{\theta, m} \phi(\theta, m | \mu^N) u_R(\theta, \sigma'(m)).$$

**Proposition 6** below shows that credible and R-IC profiles in the continuum model are approximated by credible and R-IC profiles in the  $N$ -sample model, provided  $N$  is sufficiently large. Note that in the second statement in **Proposition 6**, we distinguish a strictly credible profile  $(\lambda^*, \sigma^*)$  in the continuum model as one where  $\lambda^*$  is the unique maximizer in **Definition 1**; similarly,  $(\lambda^*, \sigma^*)$  is strictly R-IC if  $\sigma^*$  is the unique maximizer in **Definition 2**.

**Proposition 6.** 1. *Let  $(\lambda^*, \sigma^*)$  be a profile in the continuum model. If for every  $\varepsilon > 0$ , there exists a finite credible profile  $(\nu^N, \phi^N, \sigma^N)$  for some sample size  $N$ , such that  $|\nu^N - \lambda_M^*| < \varepsilon$ ,  $|\sigma^N - \sigma^*| < \varepsilon$  and  $P(|\phi^N(F_\Theta^N) - \lambda^*| < \varepsilon) > 1 - \varepsilon$ , then  $(\lambda^*, \sigma^*)$  is credible and R-IC.*

2. *Suppose  $(\lambda^*, \sigma^*)$  is a strictly credible and strictly R-IC profile in the continuum model, then for each  $\varepsilon > 0$  there exists a finite-sample credible and R-IC profile  $(\nu^N, \phi^N, \sigma^N)$*

---

<sup>23</sup>Note that our formulation of the Sender's strategy assumes that the Sender conditions her strategy only on the empirical distribution of the realized  $N$  samples, and ignores the identity of each individual sample point.

such that  $|\nu^N - \lambda_M^*| < \varepsilon$ ,  $|\sigma^N - \sigma^*| < \varepsilon$  and  $P(|\phi^N(F_\Theta^N) - \lambda^*| < \varepsilon) > 1 - \varepsilon$ .

The first statement in [Proposition 6](#) is analogous to the upper-hemicontinuity of Nash equilibrium correspondences: if a profile  $(\lambda^*, \sigma^*)$  in the continuous model can be arbitrarily approximated by credible and R-IC profiles in the finite model, then profile  $(\lambda^*, \sigma^*)$  must itself be credible and R-IC. Conversely, the second statement in [Proposition 6](#) can be interpreted in a way similar to the lower-hemicontinuity of strict Nash equilibria: if a profile  $(\lambda^*, \sigma^*)$  in the continuous model is strictly credible and strictly R-IC, then it can be arbitrarily approximated by credible and R-IC profiles in the finite model.<sup>24</sup>

### 4.3 Restriction to Pure Strategies

While our paper focuses on the Receiver playing pure strategies, the notion of credible and R-IC profiles can be extended to allow for Receiver mixing. Suppose the message space  $M$  is a Polish space that contains  $\Delta(A)$  as a subset. A profile  $(\lambda, \sigma)$  consisting of the Sender's information structure  $\lambda \in \Delta(\Theta \times M)$  and the Receiver strategy  $\sigma : M \rightarrow \Delta(A)$  is (mixed-strategy) credible if

$$\lambda \in \arg \max_{\lambda' \in D(\lambda)} \int_{\Theta \times M} \tilde{u}_S(\theta, \sigma(m)) d\lambda'(\theta, m),$$

and (mixed-strategy) R-IC if

$$\sigma \in \arg \max_{\sigma' : M \rightarrow \Delta(A)} \int_{\Theta \times M} \tilde{u}_R(\theta, \sigma'(m)) d\lambda(\theta, m),$$

where  $\tilde{u}_S : \Theta \times \Delta(A) \rightarrow \mathbb{R}$  and  $\tilde{u}_R : \Theta \times \Delta(A) \rightarrow \mathbb{R}$  are extensions of  $u_S$  and  $u_R$  to mixed strategies, respectively.

As is illustrated in the following example, allowing mixed strategies can sometimes enlarge the set of payoffs achievable through credible persuasion.<sup>25</sup> This is based on similar ideas that appeared in [Chakraborty and Harbaugh \(2010\)](#) and [Lipnowski and Ravid \(2020\)](#): by mixing actions that the Sender finds unappealing with those that she finds desirable, the Receiver can reduce the scope of the Sender's profitable deviations.

**Example 3.** Suppose  $\Theta = \{\theta_1, \theta_2\}$  with equal priors, and  $A = \{a_1, a_2, a_3\}$ . Consider the payoff matrices in [Table 4](#). In this example,  $a_3$  is the most desirable action for the Sender. We will show that without Receiver mixing, the Sender can never induce the Receiver to play  $a_3$  through credible persuasion; however, with Receiver mixing, the Sender can achieve a higher payoff by persuading the Receiver to take  $a_3$  with positive probability.

<sup>24</sup>We conjecture that for generic payoffs, an outcome distribution induced by a credible and R-IC profile can be approximated by their strict counterparts, but we have not been able to identify a proof.

<sup>25</sup>We thank a referee for suggesting this example.

| $u_S$      | $a_1$ | $a_2$ | $a_3$ |
|------------|-------|-------|-------|
| $\theta_1$ | 1     | 0     | 4     |
| $\theta_2$ | 0     | 1     | 2     |

| $u_R$      | $a_1$ | $a_2$ | $a_3$ |
|------------|-------|-------|-------|
| $\theta_1$ | 1     | 0     | -1    |
| $\theta_2$ | 0     | 1     | 1     |

Table 4: Sender and Receiver's payoffs

First, we show that without mixing, the Receiver will never play  $a_3$ . In particular, we argue that any stable outcome distribution  $\pi^*$  must satisfy  $\pi_A^*(a_3) = 0$ . Suppose by contradiction that  $\pi_A^*(a_3) > 0$ . Since  $a_3$  is weakly dominated by  $a_2$ , the Receiver will only play  $a_3$  under the point mass belief on  $\theta_2$ . It follows that  $\pi^*(\theta_2|a_3) = 1$ , so  $\pi^*(\theta_1, a_3) = 0$ . Therefore, either  $\pi(\theta_1, a_1) > 0$  or  $\pi(\theta_1, a_2) > 0$ . However, recall that  $(\theta_2, a_3)$  is in the support of  $\pi^*$  and

$$\begin{aligned} u_S(\theta_1, a_1) + u_S(\theta_2, a_3) &< u_S(\theta_1, a_3) + u_S(\theta_2, a_1), \text{ and} \\ u_S(\theta_1, a_2) + u_S(\theta_2, a_3) &< u_S(\theta_1, a_3) + u_S(\theta_2, a_2). \end{aligned}$$

So both cases violate  $u_S$ -cyclical monotonicity. This proves that only  $a_1$  and  $a_2$  can be induced in any stable outcome distribution. In fact, the best the Sender can do with credible persuasion is to fully reveal the states, which gives the Sender a payoff of 1.

Next we show that the Sender can achieve a strictly higher payoff with Receiver mixing. Consider the profile where the Sender fully reveals the state ( $\lambda(\theta_1, m_1) = \lambda(\theta_2, m_2) = \frac{1}{2}$ ), and the Receiver plays  $\sigma(m_1) = \delta_{a_1}$  and  $\sigma(m_2) = \frac{1}{2}\delta_{a_2} + \frac{1}{2}\delta_{a_3}$ , with  $\delta$  denoting the Dirac measure. This profile is clearly R-IC. Moreover, without changing the distribution of messages, the only deviation the Sender has is pairwise-swapping probability mass from  $(\theta_1, m_1)$  and  $(\theta_2, m_2)$  to be placed on  $(\theta_1, m_2)$  and  $(\theta_2, m_1)$ . This is not profitable because

$$\tilde{u}_S(\theta_1, \delta_{a_1}) + \tilde{u}_S(\theta_2, \frac{1}{2}\delta_{a_2} + \frac{1}{2}\delta_{a_3}) = 1 + 1.5 > 2 = \tilde{u}_S(\theta_1, \frac{1}{2}\delta_{a_2} + \frac{1}{2}\delta_{a_3}) + \tilde{u}_S(\theta_2, \delta_{a_1}).$$

Therefore, this strategy profile is (mixed-strategy) credible and R-IC. Moreover, the Sender achieves a strictly higher payoff of 1.25 from this mixed strategy profile than any pure-strategy credible and R-IC profile.

Despite the gap between pure and mixed strategies illustrated by the example above, some of our results can be extended to cover Receiver mixing. In [Appendix B.1](#), we provide a variant of [Theorem 1](#) ([Theorem 1\\*](#)) for the case when  $\Theta$  and  $A$  are both compact Polish spaces. If we view an outcome distribution  $\pi \in \Delta(\Theta \times \Delta(A))$  as direct recommendations for mixed strategies, [Theorem 1\\*](#) then characterizes credibility as  $\tilde{u}_S$ -cyclical monotonicity on the space  $\Theta \times \Delta(A)$ .

As a more specific example, when the Receiver's action is binary, the negative implication of [Proposition 1](#) holds even when allowing for Receiver mixing. In particular, [Proposition 1\\*](#)

in [Appendix B.1](#) extends [Proposition 1](#) to the case when  $\Theta$  and  $A$  are both compact subsets of  $\mathbb{R}$ . When the Receiver’s action is binary, the set of mixed strategies can be identified with the interval  $[0, 1]$ , and the extended payoff functions  $\tilde{u}_S$  and  $\tilde{u}_R$  preserve the super(sub-)modularity of  $u_S$  and  $u_R$ . So as a corollary of [Proposition 1\\*](#), no information can be credibly transmitted in this case, and focusing on pure strategies in [Proposition 1](#) is without loss of generality when the Receiver’s action is binary.

## 5 Conclusion

This paper offers a new notion of credibility for persuasion problems. We model a Sender who can commit to an information structure only up to the details that are observable to the Receiver. The Receiver does not observe the chosen information structure but observes the distribution of messages. This leads to a model of partial commitment where the Sender can undetectably deviate to information structures that induce the same distribution of messages. Our framework characterizes when, given the Receiver’s best response, the Sender has no profitable deviation.

We show that this consideration eliminates the prospects for credible persuasion in settings with adverse selection. In some other settings, we show that the requirement is compatible with the Sender still benefiting from persuasion. More generally, we show that our requirement translates to a cyclical monotonicity condition on the induced distribution of states and actions.

Our work also speaks to why certain industries (such as education) can effectively disclose information by utilizing their own rating systems, while some other industries (such as car dealerships) must resort to other means to address asymmetric information, such as third-party certification or warranties. Our results provide a rationale: in industries that exhibit adverse selection, the informed party cannot credibly disclose information through its own ratings even if it wishes to do so.

The notion of credibility we consider in this paper is motivated by the observability of the Sender’s message distribution. In some settings, the Receiver may observe more than the distribution of messages; for example, she may observe some further details of the information structure, such as how some states of the world map into messages. In other settings, the Receiver may observe less; e.g., she may see the average grade, but not its distribution. To capture these various cases, one would then formulate the problem of “detectable” deviations differently. We view it to be an interesting direction for future research to understand how different notions of detectability map into different conditions on the outcome distribution.

# References

- Akbarpour, Mohammad and Shengwu Li. 2020. “Credible Auctions: A Trilemma.” *Econometrica* 88 (2):425–467.
- Akerlof, George A. 1970. “The Market for ‘Lemons’: Quality Uncertainty and the Market Mechanism.” *Quarterly Journal of Economics* 84 (3):488–500.
- Aliprantis, CD and KC Border. 2006. *Infinite Dimensional Analysis: A Hitchhiker’s Guide, 3rd Edition*. Springer-Verlag.
- Alonso, Ricardo and Odilon Câmara. 2016. “Persuading Voters.” *American Economic Review* 106 (11):3590–3605.
- Beiglböck, Mathias, Martin Goldstern, Gabriel Maresch, and Walter Schachermayer. 2009. “Optimal and Better Transport Plans.” *Journal of Functional Analysis* 256 (6):1907–1927.
- Bergemann, Dirk and Stephen Morris. 2016. “Bayes Correlated Equilibrium and the Comparison of Information Structures in Games.” *Theoretical Economics* 11 (2):487–522.
- Best, James and Daniel Quigley. 2020. “Persuasion for the Long Run.” Working Paper.
- Brocas, Isabelle and Juan D Carrillo. 2007. “Influence Through Ignorance.” *The RAND Journal of Economics* 38 (4):931–947.
- Brualdi, Richard A. 2006. *Combinatorial Matrix Classes*, vol. 13. Cambridge University Press.
- Chakraborty, Archishman and Rick Harbaugh. 2007. “Comparative Cheap Talk.” *Journal of Economic Theory* 132 (1):70–94.
- . 2010. “Persuasion by Cheap Talk.” *American Economic Review* 100 (5):2361–82.
- Chang, Joseph T and David Pollard. 1997. “Conditioning as Disintegration.” *Statistica Neerlandica* 51 (3):287–317.
- Crawford, Vincent P. and Joel Sobel. 1982. “Strategic Information Transmission.” *Econometrica* 50 (6):1431–1451.
- Dworczak, Piotr and Giorgio Martini. 2019. “The Simple Economics of Optimal Persuasion.” *Journal of Political Economy* 127 (5):1993–2048.
- Frankel, Alexander. 2014. “Aligned Delegation.” *American Economic Review* 104 (1):66–83.
- Fréchette, Guillaume R, Alessandro Lizzeri, and Jacopo Perego. 2021. “Rules and commitment in communication: An experimental analysis.” *Econometrica* (Forthcoming).
- Gentzkow, Matthew and Emir Kamenica. 2017. “Bayesian persuasion with multiple senders and rich signal spaces.” *Games and Economic Behavior* 104:411–429.
- Gitmez, A. Arda and Pooya Molavi. 2022. “Polarization and Media Bias.” *arXiv preprint arXiv:2203.12698*.
- Goldstein, Itay and Yaron Leitner. 2018. “Stress Tests and Information Disclosure.” *Journal of*

- Economic Theory* 177:34–69.
- Guo, Yingni and Eran Shmaya. 2021. “Costly miscalibration.” *Theoretical Economics* 16 (2):477–506.
- Hedlund, Jonas. 2017. “Bayesian persuasion by a privately informed sender.” *Journal of Economic Theory* 167:229–268.
- Ivanov, Maxim. 2020. “Optimal Monotone Signals in Bayesian Persuasion Mechanisms.” *Economic Theory* 72:955–1000.
- Jackson, Matthew O. and Hugo F. Sonnenschein. 2007. “Overcoming Incentive Constraints by Linking Decisions.” *Econometrica* 75 (1):241–257.
- Kamenica, Emir. 2019. “Bayesian Persuasion and Information Design.” *Annual Review of Economics* 11:249–272.
- Kamenica, Emir and Matthew Gentzkow. 2011. “Bayesian Persuasion.” *American Economic Review* 101 (6):2590–2615.
- Kartik, Navin. 2009. “Strategic communication with lying costs.” *The Review of Economic Studies* 76 (4):1359–1395.
- Kartik, Navin and Weijie Zhong. 2019. “Lemonade from Lemons: Information Design and Adverse Selection.” Working Paper.
- Koessler, Frédéric and Vasiliki Skreta. 2021. “Information design by an informed designer.” Working Paper.
- Kolotilin, Anton. 2018. “Optimal Information Disclosure: A Linear Programming Approach.” *Theoretical Economics* 13 (2):607–635.
- Kolotilin, Anton and Hongyi Li. 2020. “Relational Communication.” *Theoretical Economics* (Forthcoming).
- Kuvalekar, Aditya, Elliot Lipnowski, and Joao Ramos. 2021. “Goodwill in Communication.” *Journal of Economic Theory* (Forthcoming).
- Libgober, Jonathan. 2022. “False positives and transparency.” *American Economic Journal: Microeconomics* 14 (2):478–505.
- Lipnowski, Elliot and Doron Ravid. 2020. “Cheap Talk with Transparent Motives.” *Econometrica* 88 (4):1631–1660.
- Lipnowski, Elliot, Doron Ravid, and Denis Shishkin. 2021. “Persuasion via Weak Institutions.” Working Paper.
- . 2022. “Persuasion via weak institutions.” *Journal of Political Economy* 130 (10):000–000.
- Margaria, Chiara and Alex Smolin. 2018. “Dynamic Communication with Biased Senders.” *Games and Economic Behavior* 110:330–339.
- Mas-Colell, Andreu, Michael Dennis Whinston, and Jerry R. Green. 1995. *Microeconomic Theory*,



- vol. 1. Oxford University Press New York.
- Mathevet, Laurent, David Pearce, and Ennio Stacchetti. 2022. “Reputation for A Degree of Honesty.” Working Paper.
- Matsushima, Hitoshi, Koichi Miyazaki, and Nobuyuki Yagi. 2010. “Role of Linking Mechanisms in Multitask Agency with Hidden Information.” *Journal of Economic Theory* 145 (6):2241–2259.
- Meng, Delong. 2021. “On the Value of Repetition for Communication Games.” *Games and Economic Behavior* 127:227–246.
- Mensch, Jeffrey. 2021. “Monotone Persuasion.” *Games and Economic Behavior* (Forthcoming).
- Min, Daehong. 2021. “Bayesian Persuasion under Partial Commitment.” *Economic Theory* 72:743–764.
- Nguyen, Anh and Teck Yong Tan. 2021. “Bayesian Persuasion with Costly Messages.” *Journal of Economic Theory* 193:105212.
- Ostrovsky, Michael and Michael Schwarz. 2010. “Information Disclosure and Unraveling in Matching Markets.” *American Economic Journal: Microeconomics* 2 (2):34–63.
- Pei, Harry. 2020. “Repeated Communication with Private Lying Cost.” Working Paper.
- Perez-Richet, Eduardo. 2014. “Interim bayesian persuasion: First steps.” *American Economic Review* 104 (5):469–474.
- Perez-Richet, Eduardo and Vasiliki Skreta. 2021. “Test Design under Falsification.” Working Paper.
- Rahman, David. 2010. “Detecting Profitable Deviations.” Working Paper.
- Rayo, Luis and Ilya Segal. 2010. “Optimal Information Disclosure.” *Journal of Political Economy* 118 (5):949–987.
- Renault, Jérôme, Eilon Solan, and Nicolas Vieille. 2013. “Dynamic Sender–Receiver games.” *Journal of Economic Theory* 148 (2):502–534.
- Rochet, Jean-Charles. 1987. “A Necessary and Sufficient Condition for Rationalizability in a Quasi-linear Context.” *Journal of Mathematical Economics* 16 (2):191–200.
- Santambrogio, Filippo. 2015. “Optimal Transport for Applied Mathematicians.” *Birkhäuser, NY* 55 (58-63):94.
- Taneva, Ina. 2019. “Information Design.” *American Economic Journal: Microeconomics* 11 (4):151–85.
- Topkis, Donald M. 2011. *Supermodularity and Complementarity*. Princeton University Press.
- Villani, Cédric. 2008. *Optimal Transport: Old and New*, vol. 338. Springer Science & Business Media.
- Zapechelnyuk, Andriy. 2023. “On the equivalence of information design by uninformed and informed principals.” *Economic Theory* :1–17.

# A Appendix

## A.1 Proof of Theorem 1

The following lemma, which will play a key role in the proof of Theorem 1, is a finite version of Theorem 5.10 of Villani (2008). Below we present a direct proof of the lemma to better illustrate the intuition behind Theorem 1.

**Lemma 2.** *Suppose both  $X$  and  $Y$  are finite sets, and  $u : X \times Y \rightarrow \mathbb{R}$  is a real function. Let  $\mu \in \Delta(X)$  and  $\nu \in \Delta(Y)$  be two probability measure on  $X$  and  $Y$  respectively, and  $\Pi(\mu, \nu)$  be the set of joint probability measure on  $X \times Y$  such that the marginals on  $X$  and  $Y$  are  $\mu$  and  $\nu$ . The following two statements are equivalent:*

1.  $\pi^* \in \arg \max_{\pi \in \Pi(\mu, \nu)} \sum_{x, y} \pi(x, y) u(x, y);$
2.  $\pi^*$  is  $u$ -cyclically monotone. That is, for any  $n$  and  $(x_1, y_1), \dots, (x_n, y_n) \in \text{supp}(\pi^*),$

$$\sum_{i=1}^n u(x_i, y_i) \geq \sum_{i=1}^n u(x_i, y_{i+1})$$

where  $y_{n+1} \equiv y_1$ .

*Proof.* ( $1 \Rightarrow 2$ ) To see the necessity of  $u$ -cyclical monotonicity, suppose by contraposition that  $\pi^*$  is not  $u$ -cyclically monotone, which implies that there exists a sequence  $(x_1, y_1), \dots, (x_n, y_n) \in \text{supp}(\pi^*)$  such that

$$\sum_{i=1}^n u(x_i, y_i) < \sum_{i=1}^n u(x_i, y_{i+1}),$$

where  $y_{n+1} = y_1$ . Take  $0 < \varepsilon \leq \frac{1}{n} \min_{i=1, \dots, n} \pi^*(x_i, y_i),$ <sup>26</sup> and define

$$\pi^\varepsilon \equiv \pi^* + \varepsilon \sum_{i=1}^n [\delta_{(x_i, y_{i+1})} - \delta_{(x_i, y_i)}],$$

---

<sup>26</sup>The scaling factor  $\frac{1}{n}$  is added to ensure that  $\pi^\varepsilon$  is a positive measure, in case the same pair  $(x, y)$  appears multiple times in the sequence.

where  $\delta_{(x,y)}$  denotes the Dirac measure on  $(x,y)$ . Note that  $\pi^\varepsilon \in \Pi(\mu, \nu)$  and satisfies

$$\begin{aligned} & \sum_{x,y} u(x,y) \pi^\varepsilon(x,y) \\ &= \sum_{x,y} u(x,y) \pi^*(x,y) + \varepsilon \left[ \sum_{i=1}^n u(x_i, y_{i+1}) - \sum_{i=1}^n u(x_i, y_i) \right] \\ &> \sum_{x,y} u(x,y) \pi^*(x,y), \end{aligned}$$

which implies  $\pi^* \notin \arg \max_{\pi \in \Pi(\mu, \nu)} \sum_{x,y} \pi(x,y) u(x,y)$ .

(1  $\Leftarrow$  2) First note that  $\pi^*$  being  $u$ -cyclically monotone is equivalent to: for any  $n$  and  $(x_1, y_1), \dots, (x_n, y_n) \in \text{supp}(\pi^*)$ , and for any permutation  $s : \{1, \dots, n\} \rightarrow \{1, \dots, n\}$ ,

$$\sum_{i=1}^n u(x_i, y_i) \geq \sum_{i=1}^n u(x_i, y_{s(i)}).$$

This is because any permutation can be written as the composition of disjoint cycles.

We now prove the sufficiency of  $u$ -cyclical monotonicity by contraposition. Suppose there exists  $\pi' \in \Pi(\mu, \nu)$  such that

$$\sum_{x,y} \pi'(x,y) u(x,y) > \sum_{x,y} \pi^*(x,y) u(x,y). \quad (3)$$

We will show that there exists a sequence  $(x_1, y_1), \dots, (x_n, y_n) \in \text{supp}(\pi^*)$  and a permutation  $s : \{1, \dots, n\} \rightarrow \{1, \dots, n\}$ , such that

$$\sum_{i=1}^n u(x_i, y_{s(i)}) > \sum_{i=1}^n u(x_i, y_i).$$

The argument proceeds in three steps.

Step 1: Approximate  $\pi^*$  and  $\pi'$  with  $\tilde{\pi}^*$  and  $\tilde{\pi}'$ , respectively, while preserving inequality (3): both  $\tilde{\pi}^*$  and  $\tilde{\pi}'$  are joint distributions that share the same rational marginals; in addition,  $\tilde{\pi}^*$  shares the same support as  $\pi^*$ .

Choose  $\varepsilon_0 > 0$  so that

$$\sum_{x,y} \pi^*(x,y) u(x,y) < \sum_{x,y} \pi'(x,y) u(x,y) - \varepsilon_0$$

By continuity, there exists  $\delta_1 > 0$  such that for all  $|\pi - \pi^*| < \delta_1$  we have

$$\sum_{x,y} \pi(x,y)u(x,y) < \sum_{x,y} \pi'(x,y)u(x,y) - \varepsilon_0/2 \quad (4)$$

By Lemma 10, there exists  $\delta_2 > 0$  such that for all  $|\tilde{\mu} - \mu| < \delta_2$  and  $|\tilde{\nu} - \nu| < \delta_2$ , there exists  $\pi \in \Pi(\tilde{\mu}, \tilde{\nu})$  with

$$\sum_{x,y} \pi(x,y)u(x,y) > \sum_{x,y} \pi'(x,y)u(x,y) - \varepsilon_0/2. \quad (5)$$

Let  $\delta_3 \equiv \min_{x,y} \{\pi^*(x,y) : \pi^*(x,y) > 0\}$  denote the smallest probability weight among the the support of  $\pi^*$ .

Now let  $\delta = \min\{\delta_1, \frac{\delta_2}{|X| \times |Y|}, \delta_3\}$  and consider a rational joint distribution  $\tilde{\pi}^* \in \mathbb{Q}^{X \times Y} \cap \Delta(X \times Y)$  such that  $|\tilde{\pi}^* - \pi^*| < \delta$ . Note that  $\text{supp}(\tilde{\pi}^*) = \text{supp}(\pi^*)$ . By inequality (4),

$$\sum_{x,y} \tilde{\pi}^*(x,y)u(x,y) < \sum_{x,y} \pi'(x,y)u(x,y) - \varepsilon_0/2.$$

Furthermore, the marginals of  $\tilde{\pi}^*$ ,  $p \equiv \tilde{\pi}_X^*$  and  $q \equiv \tilde{\pi}_Y^*$ , are also rational and satisfy  $|p - \mu| < \delta_2$  and  $|q - \nu| < \delta_2$ . By inequality (5) there exists  $\tilde{\pi}' \in \Pi(p, q)$  such that

$$\sum_{x,y} \tilde{\pi}'(x,y)u(x,y) > \sum_{x,y} \pi'(x,y)u(x,y) - \varepsilon_0/2,$$

so

$$\sum_{x,y} \tilde{\pi}'(x,y)u(x,y) > \sum_{x,y} \tilde{\pi}^*(x,y)u(x,y). \quad (6)$$

*Step 2: Normalize and transform the above two joint distributions with the same rational marginals,  $\tilde{\pi}^*$  and  $\tilde{\pi}'$ , into doubly stochastic matrices. Through the Birkhoff-von Neumann theorem, express inequality (6) in terms of permutation matrices.*

Let  $N$  be an integer such that  $Np(x)$  and  $Nq(y)$  are integers for all  $x \in X$  and  $y \in Y$ . Let  $S : X \rightarrow 2^{\{1, \dots, N\}}$  be a partition of  $\{1, \dots, N\}$  such that  $|S(x)| = Np(x)$  for each  $x \in X$ ; similarly, let  $T : Y \rightarrow 2^{\{1, \dots, N\}}$  be a partition of  $\{1, \dots, N\}$  such that  $|T(y)| = Nq(y)$  for each  $y \in Y$ . For each  $i = 1, \dots, N$ , let  $\tilde{x}_i \equiv \{x : i \in S(x)\}$  denote the  $x \in X$  indexing the partition that contains  $i$ ; similarly, for each column  $j$  let  $\tilde{y}_j \equiv \{y : j \in T(y)\}$  denote the  $y \in Y$  indexing the partition that contains  $j$ .

Consider the matrix  $[B_{ij}^*]_{1 \leq i, j \leq N}$  defined by

$$B_{ij}^* = \frac{\tilde{\pi}^*(\tilde{x}_i, \tilde{y}_j)}{Np(\tilde{x}_i)q(\tilde{y}_j)} \quad \text{for all } 1 \leq i, j \leq N.$$

And the matrix  $[B'_{ij}]_{1 \leq i, j \leq N}$  defined by

$$B'_{ij} = \frac{\tilde{\pi}'(\tilde{x}_i, \tilde{y}_j)}{Np(\tilde{x}_i)q(\tilde{y}_j)} \quad \text{for all } 1 \leq i, j \leq N.$$

Notice that the matrix  $B^*$  is doubly stochastic: for any  $i$  we have

$$\sum_j B^*_{ij} = \sum_{y \in Y} \left( \frac{\tilde{\pi}^*(\tilde{x}_i, y)}{Np(\tilde{x}_i)q(y)} \cdot Nq(y) \right) = \frac{p(\tilde{x}_i)}{p(\tilde{x}_i)} = 1.$$

Similarly, we can show that  $\sum_i B^*_{ij} = 1$  for every  $j$ . And following similar arguments, the matrix  $B'$  is also doubly stochastic.

Note that

$$\begin{aligned} \sum_{i,j} B^*_{ij} \cdot u(\tilde{x}_i, \tilde{y}_j) &= \sum_{i,j} \frac{\tilde{\pi}^*(\tilde{x}_i, \tilde{y}_j)}{Np(\tilde{x}_i)q(\tilde{y}_j)} u(\tilde{x}_i, \tilde{y}_j) \\ &= \sum_{x,y} \frac{\tilde{\pi}^*(x, y)}{Np(x)q(y)} \cdot Np(x) \cdot Nq(y) \cdot u(x, y) = N \sum_{x,y} \tilde{\pi}^*(x, y) u(x, y), \end{aligned}$$

and similarly,

$$\sum_{i,j} B'_{ij} \cdot u(\tilde{x}_i, \tilde{y}_j) = N \sum_{x,y} \tilde{\pi}'(x, y) u(x, y).$$

Now since

$$\sum_{x,y} \tilde{\pi}'(x, y) u(x, y) > \sum_{x,y} \tilde{\pi}^*(x, y) u(x, y),$$

we have

$$\sum_{i,j} B'_{ij} \cdot u(\tilde{x}_i, \tilde{y}_j) > \sum_{i,j} B^*_{ij} \cdot u(\tilde{x}_i, \tilde{y}_j).$$

Let  $\mathcal{P}$  denote the set of  $N \times N$  permutation matrices. By the Birkhoff-von Neumann theorem, both  $B^*$  and  $B'$  are in the convex hull of  $\mathcal{P}$ . It follows that there exist permutation matrices  $P^*$  and  $P'$  such that

$$\sum_{i,j} P'_{ij} \cdot u(\tilde{x}_i, \tilde{y}_j) > \sum_{i,j} P^*_{ij} \cdot u(\tilde{x}_i, \tilde{y}_j), \quad (7)$$

and in addition,  $P^*_{ij} = 1$  implies that the corresponding entry in  $B^*$  satisfies  $B^*_{ij} > 0$ .

*Step 3: Convert inequality (7) into a cyclical deviation.*

Note that the permutation matrix  $P^*$  is equivalent to a mapping  $t : \{1, \dots, N\} \rightarrow \{1, \dots, N\}$

such that  $P_{ij}^* = 1$  if and only if  $j = t(i)$ . So

$$\sum_{i,j} P_{ij}^* \cdot u(\tilde{x}_i, \tilde{y}_j) = \sum_i u(\tilde{x}_i, \tilde{y}_{t(i)}).$$

In particular, every element of  $\{(\tilde{x}_i, \tilde{y}_{t(i)})\}_{i=1}^N$  is in the support of  $\tilde{\pi}^*$ , since  $P_{ij}^* = 1$  implies  $B_{ij}^* > 0$ , which further implies  $\tilde{\pi}^*(\tilde{x}, \tilde{y}_{t(i)}) > 0$ . Since  $\pi^*$  and  $\tilde{\pi}^*$  share the same support, every element of  $\{(\tilde{x}_i, \tilde{y}_{t(i)})\}_{i=1}^N$  is in the support of  $\pi^*$  as well.

Let  $t' : \{1, \dots, N\} \rightarrow \{1, \dots, N\}$  denote the permutation mapping induced by the matrix  $P'_{ij}$ , so

$$\sum_{i,j} P'_{ij} \cdot u(\tilde{x}_i, \tilde{y}_j) = \sum_i u(\tilde{x}_i, \tilde{y}_{t'(i)}).$$

It follows that inequality (7) can be re-written as

$$\sum_i u(\tilde{x}_i, \tilde{y}_{t'(i)}) > \sum_i u(\tilde{x}_i, \tilde{y}_{t(i)}). \quad (8)$$

Since  $t$  and  $t'$  are both permutations, there exists a permutation  $s : \{1, \dots, N\} \rightarrow \{1, \dots, N\}$  such that  $s(t(i)) = t'(i)$ . Consider the sequence  $\{(x_i, y_i)\}_{i=1}^N$  defined by

$$x_i = \tilde{x}_i \text{ and } y_i = \tilde{y}_{t(i)} \text{ for } i = 1, \dots, N.$$

Then (8) becomes

$$\sum_{i=1}^N u(x_i, y_i) < \sum_{i=1}^N u(x_i, y_{s(i)})$$

where  $(x_i, y_i) \in \text{supp}(\pi^*)$  for each  $1 \leq i \leq N$ , which violates  $u$ -cyclical monotonicity.  $\square$

*Proof of Theorem 1.* For the “if” direction, suppose  $\pi$  is  $u_R$ -obedient,  $u_S$ -cyclically monotone, and satisfies  $\pi_\Theta = \mu_0$ . The proof is by construction.

Since  $\pi_\Theta = \mu_0$ , we can construct an information structure  $(M, \lambda^*)$  by setting  $M = A$  and  $\lambda^* = \pi$ ; furthermore, let  $\sigma^*$  be the identity map from  $M$  to  $A$ . It is straightforward to see that the profile  $(\lambda^*, \sigma^*)$  induces the outcome distribution  $\pi$ . We first show that  $(\lambda^*, \sigma^*)$  is R-IC. Since  $\pi$  is  $u_R$ -obedient, we have that for each  $a \in A$ ,

$$a \in \arg \max_{a'} \sum_{\Theta} u_R(\theta, a') \pi(\theta, a).$$

Since  $\sigma^*$  is an identity map, it follows that for each  $m \in M$ ,

$$\sigma^*(m) \in \arg \max_{a'} \sum_{\Theta} u_R(\theta, a') \pi(\theta, \sigma^*(m)).$$

Furthermore, since  $\lambda^* = \pi$  and  $\sigma^*$  is injective, we have  $\lambda^*(\theta, m) = \pi(\theta, \sigma^*(m))$  for all  $\theta \in \Theta$  and  $m \in M$ . So

$$\sigma^* \in \arg \max_{\sigma: M \rightarrow A} \sum_{\Theta \times M} u_R(\theta, \sigma(m)) \lambda^*(\theta, m),$$

which means  $\sigma^*$  is a best response to  $\lambda^*$ .

It remains to show that the Sender does not benefit from choosing any other information structure in  $D(\lambda^*)$ . Observe that since  $\pi$  is  $u_S$ -cyclically monotone, every sequence  $(\theta_1, a_1), \dots, (\theta_n, a_n)$  in  $\text{supp}(\pi)$  where  $a_{n+1} \equiv a_1$  satisfies

$$\sum_{i=1}^n u_S(\theta_i, a_i) \geq \sum_{i=1}^n u_S(\theta_i, a_{i+1}).$$

Since  $\lambda^* = \pi$  and  $\sigma^*$  is the identity mapping, this further implies

$$\sum_{i=1}^n u_S(\theta_i, \sigma^*(m_i)) \geq \sum_{i=1}^n u_S(\theta_i, \sigma^*(m_{i+1}));$$

for every sequence  $(\theta_1, m_1), \dots, (\theta_n, m_n) \in \text{supp}(\lambda^*)$  with  $m_{n+1} = m_1$ . In addition,  $\lambda_\theta^* = \mu_0$  and  $\lambda_M^* = \lambda_M^*$  by construction. By [Lemma 2](#),  $\lambda^*$  satisfies

$$\lambda^* \in \arg \max_{\lambda \in D(\lambda^*)} \sum_{\Theta \times M} u_S(\theta, \sigma(m)) \lambda(\theta, m)$$

which means  $\lambda^*$  is Sender optimal conditional on its message distribution.

For the “only if” direction, suppose  $\pi$  is stable and thus induced by a credible and R-IC profile  $(\lambda^*, \sigma^*)$ . Since  $\sigma^*$  best responds to the messages from  $\lambda^*$ , the  $u_R$ -obedience of  $\pi$  follows from [Bergemann and Morris \(2016\)](#).

It remains to show that  $\pi$  is  $u_S$ -cyclical monotone. Suppose by contradiction that  $\pi$  is not  $u_S$ -cyclically monotone, which implies that there exists a sequence  $(\theta_1, a_1), \dots, (\theta_n, a_n) \in \text{supp}(\pi)$  such that

$$\sum_{i=1}^n u_S(\theta_i, a_i) < \sum_{i=1}^n u_S(\theta_i, a_{i+1}),$$

where  $a_{n+1} = a_1$ . Since  $\pi$  is induced by  $(\lambda^*, \sigma^*)$ , for each  $i = 1, \dots, n$  there exists  $m_i$  such that  $m_i \in \sigma^{*-1}(a_i)$  and  $(\theta_i, m_i) \in \text{supp}(\lambda^*)$ , so we have a sequence  $(\theta_1, m_1), \dots, (\theta_n, m_n) \in \text{supp}(\lambda^*)$  that satisfies

$$\sum_{i=1}^n u_S(\theta_i, \sigma^*(m_i)) < \sum_{i=1}^n u_S(\theta_i, \sigma^*(m_{i+1})), \quad (9)$$

where  $m_{n+1} = m_1$ . Define  $v(\theta, m) \equiv u_S(\theta, \sigma^*(m))$ . Since  $(\lambda^*, \sigma^*)$  is credible, we have

$$\lambda^* \in \arg \max_{\lambda \in D(\lambda^*)} \sum_{\Theta \times M} v(\theta, m) \lambda(\theta, m).$$

**Lemma 2** implies that  $\lambda^*$  is  $v$ -cyclically monotone. Since  $(\theta_1, m_1), \dots, (\theta_n, m_n)$  is in  $\text{supp}(\lambda^*)$ , the  $v$ -cyclical monotonicity of  $\lambda^*$  implies

$$\sum_{i=1}^n u_S(\theta_i, \sigma^*(m_i)) \geq \sum_{i=1}^n u_S(\theta_i, \sigma^*(m_{i+1}))$$

where  $m_{n+1} = m_1$ , which is a contradiction to (9). So  $\pi$  must be  $u_S$ -cyclically monotone.  $\square$

## A.2 Proof of Corollary 1

The Sender-optimal stable outcome distribution is the solution to the following problem:

$$\begin{aligned} & \max_{\pi \in \Delta(\Theta \times A)} \sum_{\theta, a} \pi(\theta, a) u_S(\theta, a) \\ \text{s.t. } & \sum_{\theta} \pi(\theta|a) u_R(\theta, a) \geq \sum_{\theta} \pi(\theta|a) u_R(\theta, a') \text{ for all } a \in \text{supp}(\pi_A) \text{ and } a' \in A, \\ & \pi \text{ is } u_S\text{-cyclically monotone.} \end{aligned}$$

We first argue that the feasible region in the optimization program above is compact, so there exists a Sender-optimal stable outcome distribution. To this end, let  $\Pi_C \equiv \{\pi \in \Delta(\Theta \times A) : \pi \text{ is } u_S\text{-cyclically monotone}\}$  denote the set of  $u_S$ -cyclically monotone outcome distributions. It suffices to argue that  $\Pi_C$  is compact.

Let  $\mathcal{O} \equiv \{\text{supp}(\pi) : \pi \in \Pi_C\}$  denote the set of the supports of the distributions in  $\Pi_C$ . For each such support  $O \in \mathcal{O}$ , let  $\Pi_O \equiv \{\pi \in \Delta(\Theta \times A) : \text{supp}(\pi) \subseteq O\}$  denote the set of outcome distributions whose support is contained within  $O$ . Note that every distribution in the set  $\Pi_O$  is  $u_S$ -cyclically monotone, so we have  $\Pi_C = \cup_{O \in \mathcal{O}} \Pi_O$ . In addition, for each  $O \in \mathcal{O}$ , the set  $\Pi_O$  is closed since it is defined by equality constraints:  $\Pi_O = \{\pi \in \Delta(\Theta \times A) : \pi(\theta, a) = 0 \ \forall (\theta, a) \notin O\}$ . The set  $\Pi_C = \cup_{O \in \mathcal{O}} \Pi_O$  is a finite union of closed and bounded sets, and is therefore compact, so there exists a Sender-optimal stable outcome distribution, which we shall denote by  $\pi^*$ .

Let  $v^* \equiv \sum_a \pi_A^*(a) \sum_{\theta} \pi^*(\theta|a) u_S(\theta, a)$  denote the Sender's value from  $\pi^*$ , and  $A^* \equiv$



$\text{supp}(\pi_A^*)$  denote the support of  $\pi^*$ 's action distribution. Consider the following set

$$\mathcal{E} \equiv \left\{ \left( \pi^*(\cdot|a), \sum_{\theta} u_S(\theta, a) \pi^*(\theta|a) \right) \in \mathbb{R}^{|\Theta|} \mid a \in A^* \right\},$$

where  $\pi^*(\cdot|a) \in \Delta(\Theta)$  and  $\sum_{\theta} u_S(\theta, a) \pi^*(\theta|a) \in \mathbb{R}$  denote the posterior belief and Sender's value conditioning on  $a$ , respectively. Let  $(\mu_a, v_a)$  denote a generic element of  $\mathcal{E}$ .

Recall that  $\mu_0 \in \Delta(\Theta)$  is the prior distribution over states. Clearly  $(\mu_0, v^*) \in \text{conv}(\mathcal{E}) \subset \mathbb{R}^{|\Theta|}$ . We will show that  $(\mu_0, v^*)$  can be represented as a convex combination of at most  $|\Theta|$  number of points in  $\mathcal{E}$ . First we show that  $(\mu_0, v^*)$  must be a boundary point of  $\text{conv}(\mathcal{E})$ . Suppose not, then there exists  $\hat{p} \in \Delta(A^*)$  and  $\{(\mu_a, v_a)\}_{a \in A^*} \in \mathcal{E}$  such that  $\sum_{a \in A^*} (\mu_a, v_a) \hat{p}(a) = (\pi_0, \hat{v})$  where  $\hat{v} > v^*$ . Let  $\hat{\pi} \in \Delta(\Theta \times A)$  be the outcome distribution induced by  $\hat{p}$ : that is,

$$\hat{\pi}(\theta, a) = \begin{cases} \hat{p}(a) \pi^*(\theta|a) & \text{for all } a \in A^* \text{ and } \theta \in \Theta, \\ 0 & \text{otherwise.} \end{cases}$$

Then  $\hat{\pi}$  satisfies both obedience and cyclical monotonicity because the support of  $\hat{\pi}$  is a subset of  $\pi^*$ . Moreover, it yields a strictly higher value to the Sender, which contradicts  $\pi^*$  being the Sender-optimal stable outcome distribution. Therefore,  $(\mu_0, v^*)$  is on the boundary of  $\text{conv}(\mathcal{E})$ , and thus on a face of  $\text{conv}(\mathcal{E})$ , which has  $|\Theta| - 1$  dimensions. By Carathéodory's Theorem,  $(\mu_0, v^*)$  can be represented as a convex combination of at most  $|\Theta|$  number of points in  $\mathcal{E}$  according to some mixture probabilities  $\tilde{p} \in \Delta(A^*)$ .

Let  $\tilde{\pi}$  be the outcome distribution obtained from the same construction as that for  $\hat{\pi}$  above but with  $\tilde{p}$  replacing  $\hat{p}$ . By construction,  $|\text{supp}(\tilde{\pi}_A)| \leq \min\{|\Theta|, |A|\}$ , and the Sender's value from  $\tilde{p}$  is also  $v^*$ . In addition,  $\tilde{\pi}$  satisfies both obedience and cyclical monotonicity because the support of  $\tilde{\pi}$  is a subset of  $\pi^*$ , so  $\tilde{\pi}$  is a Sender-optimal stable outcome distribution. Using the construction outlined in the “only if” direction of the proof of [Theorem 1](#), we can construct a “direct recommendation” Sender-optimal credible and R-IC profile  $(\lambda^*, \sigma^*)$ , where  $\lambda^*$  has no more than  $\min\{|\Theta|, |A|\}$  messages.

### A.3 Proof of [Corollary 2](#)

*Proof.* We prove that if there exists a sequence  $(\theta_1, a_1), \dots, (\theta_n, a_n) \in \text{supp}(\pi)$  with  $n > \min\{|\Theta|, |A|\}$  such that

$$u_S(\theta_1, a_1) + \dots + u_S(\theta_n, a_n) < u_S(\theta_1, a_2) + \dots + u_S(\theta_n, a_1),$$

then there exists a sequence  $(\theta_1, a_1), \dots, (\theta_m, a_m) \in \text{supp}(\pi)$  with  $m < n$  such that

$$u_S(\theta_1, a_1) + \dots + u_S(\theta_m, a_m) < u_S(\theta_1, a_2) + \dots + u_S(\theta_m, a_1).$$

Suppose by contradiction that there exists a sequence  $(\theta_1, a_1), \dots, (\theta_n, a_n)$  with  $n > \min\{|\Theta|, |A|\}$  such that

$$u_S(\theta_1, a_1) + \dots + u_S(\theta_n, a_n) < u_S(\theta_1, a_2) + \dots + u_S(\theta_n, a_1).$$

and that for all sequences with length  $m < n$ ,

$$u_S(\theta_1, a_1) + \dots + u_S(\theta_m, a_m) \geq u_S(\theta_1, a_2) + \dots + u_S(\theta_m, a_1). \quad (10)$$

Suppose  $\min\{|\Theta|, |A|\} = |A|$  (a similar argument works for  $\min\{|\Theta|, |A|\} = |\Theta|$ ), then there exists  $a^*$  that appears twice in the sequence  $(\theta_1, a_1), \dots, (\theta_n, a_n)$ . Without loss let  $a_1 = a_k = a^*$  with  $1 < k \leq n$ . Then

$$\begin{aligned} u_S(\theta_1, a_1) + \dots + u_S(\theta_n, a_n) &= [u_S(\theta_1, a_1) + \dots + u_S(\theta_{k-1}, a_{k-1})] + [u_S(\theta_k, a_k) + \dots + u_S(\theta_n, a_n)] \\ &\geq [u_S(\theta_1, a_2) + \dots + u_S(\theta_{k-1}, a_1)] + [u_S(\theta_k, a_{k+1}) + \dots + u_S(\theta_n, a_k)] \\ &= [u_S(\theta_1, a_2) + \dots + u_S(\theta_{k-1}, a_k)] + [u_S(\theta_k, a_{k+1}) + \dots + u_S(\theta_n, a_{n+1})] \end{aligned}$$

where the inequality follows from (10), and the second equality holds because  $a_k = a_1$ . This is a contradiction.  $\square$

## A.4 Proof of Lemma 1

In light of Remark 1, we shall prove Lemma 1 without assuming that the order on either  $\Theta$  or  $A$  is antisymmetric. Suppose  $(\Theta, \succeq_1)$  and  $(A, \succeq_2)$  are finite ordered sets, and  $\succeq_1, \succeq_2$  are weak orders (complete and transitive). The notions of supermodularity and comonotonicity are extended naturally with weak orders  $\succeq_1$  and  $\succeq_2$  replacing the total orders on  $\Theta$  and  $A$ .

In particular, we say a function  $u : \Theta \times A \rightarrow \mathbb{R}$  is supermodular if for any  $\theta \succeq \theta'$  and  $a \succeq a'$  we have

$$u(\theta, a) + u(\theta', a') \geq u(\theta, a') + u(\theta', a);$$

the function is strictly supermodular if in addition for any  $\theta \succ \theta'$  and  $a \succ a'$ ,

$$u(\theta, a) + u(\theta', a') > u(\theta, a') + u(\theta', a).$$

An outcome distribution  $\pi$  is comonotone if for any  $(\theta, a), (\theta', a') \in \text{supp}(\pi)$ ,  $\theta \succ \theta'$  implies

$a \succeq a'$ . We shall prove the following result, which is a restatement of [Lemma 1](#) but based on the weak orders  $\succeq_1$  and  $\succeq_2$ .

**Lemma 1\*.** If  $u_S$  is supermodular, then every comonotone outcome distribution is  $u_S$ -cyclically monotone. Furthermore, if  $u_S$  is strictly supermodular, then every  $u_S$ -cyclically monotone outcome distribution is also comonotone.

We begin the proof by establishing the following lemma.

**Lemma 3.** Let  $t : \{1, \dots, n\} \rightarrow \{1, \dots, n\}$  be a bijection. Suppose  $t$  is not the identity mapping, then there exists  $k^*$  such that  $t(k^*) > k^*$  and  $t(t(k^*)) < t(k^*)$ .

*Proof.* Define  $K := \{k \in \{1, \dots, n\} : t(k) \neq k\}$ . Since  $t$  is not the identity mapping,  $K$  is nonempty. Since  $t$  is a bijection,  $t(k) \neq k$  if and only if  $t(t(k)) \neq t(k)$ , so  $K$  is  $t$ -invariant. Let  $k^* = t^{-1}(\max K) \in K$ , then  $k^* < \max K = t(k^*)$  and  $t(k^*) = \max K > t(\max K) = t(t(k^*))$ .  $\square$

*Proof of Lemma 1\*.* First, we show that comonotonicity implies  $u_S$ -cyclical monotonicity when  $u_S$  is supermodular. Suppose an outcome distribution  $\pi \in \Delta(\Theta \times A)$  is comonotone, then the product order of  $\succeq_1$  and  $\succeq_2$  is also a weak order on  $\text{supp}(\pi)$ . Take any sequence  $(\theta_1, a_1), \dots, (\theta_n, a_n) \in \text{supp}(\pi)$  and assume without loss of generality that  $(\theta_i, a_i)$  is non-decreasing in  $i \in \{1, \dots, n\}$  with respect to the product order. We will show that for any permutation  $t : \{1, \dots, n\} \rightarrow \{1, \dots, n\}$ ,

$$u_S(\theta_1, a_1) + \dots + u_S(\theta_n, a_n) \geq u_S(\theta_1, a_{t(1)}) + \dots + u_S(\theta_n, a_{t(n)}),$$

which then proves the statement. In particular, for each permutation  $t$ , let  $v(t) \equiv u_S(\theta_1, a_{t(1)}) + \dots + u_S(\theta_n, a_{t(n)})$  denote the value obtained from summing  $u_S$  according to the state-action pairings in  $t$  and let  $I$  denote the identity map. We show that  $v(I) \geq v(t)$  for every permutation  $t$ .

To this end, take any permutation  $t$  that is not an identity mapping, and let  $l(t)$  denote the number of fixed points of  $t$  (which may be zero). By [Lemma 3](#), there exists  $k^*$  such that  $t(k^*) > k^*$  and  $t(t(k^*)) < t(k^*)$ . The supermodularity of  $u_S$  implies

$$u_S(\theta_{t(k^*)}, a_{t(k^*)}) + u_S(\theta_{k^*}, a_{t(t(k^*))}) \geq u_S(\theta_{k^*}, a_{t(k^*)}) + u_S(\theta_{t(k^*)}, a_{t(t(k^*))}). \quad (11)$$

Define a new permutation  $\hat{t}$  so that  $k$  is mapped to  $t(t(k))$  while  $t(k)$  is mapped to  $t(k)$ , while

all other pairings remain unchanged. Formally,

$$\hat{t}(k) = \begin{cases} t(k) & \text{for all } k \neq k^*, t(k^*) \\ t(t(k^*)) & \text{if } k = k^* \\ t(k^*) & \text{if } k = t(k^*) \end{cases}$$

By (11), we have

$$u_S(\theta_1, a_{\hat{t}(1)}) + \dots + u_S(\theta_n, a_{\hat{t}(n)}) \geq u_S(\theta_1, a_{t(1)}) + \dots + u_S(\theta_n, a_{t(n)}),$$

so we have constructed another permutation  $\hat{t}$  with  $v(\hat{t}) \geq v(t)$  and  $l(\hat{t}) = l(t) + 1$ . Each time we iterate the process above,  $v(\cdot)$  weakly increases while the number of fixed points increases by one. Since  $n < \infty$ , the iteration terminates at the identity map  $I$ , so  $v(I) \geq v(t)$  for every permutation  $t$ .

Next, suppose  $u_S$  is strictly supermodular. We want to show that  $u_S$ -cyclical monotonicity implies comonotonicity. We prove this statement by contraposition: suppose that an outcome distribution  $\pi$  is not comonotone, we will show that  $\pi$  is not  $u_S$ -cyclically monotone. Since  $\pi$  is not comonotone, there exists  $(\theta, a), (\theta', a') \in \text{supp}(\pi)$  such that  $\theta \prec \theta', a \succ a'$ . Since  $u_S$  is strictly supermodular, we have

$$u_S(\theta, a) + u_S(\theta', a') < u_S(\theta, a') + u_S(\theta', a). \quad (12)$$

Consider a cycle of length 2 where  $(\theta_1, a_1) = (\theta, a)$  and  $(\theta_2, a_2) = (\theta', a')$ , then inequality (12) above implies that  $\pi$  is not  $u_S$ -cyclically monotone.  $\square$

## A.5 Proof of Proposition 1

Let  $\pi$  be a stable outcome distribution, and suppose by contradiction that there exist two distinct actions  $a_1, a_2 \in \text{supp}(\pi_A)$ , say  $a_1 < a_2$ . Let  $I_1 \equiv \{\theta \in \Theta \mid \pi(\theta, a_1) > 0\}$  and  $I_2 \equiv \{\theta \in \Theta \mid \pi(\theta, a_2) > 0\}$  be the states associated with  $a_1$  and  $a_2$  in the support of  $\pi$ , respectively. By Theorem 1, since  $\pi$  is stable, it must be  $u_R$ -obedient, which implies

$$\sum_{\theta \in I_1} [u_R(\theta, a_1) - u_R(\theta, a_2)] \frac{\pi(\theta, a_1)}{\pi_A(a_1)} \geq 0 \geq \sum_{\theta' \in I_2} [u_R(\theta', a_1) - u_R(\theta', a_2)] \frac{\pi(\theta', a_2)}{\pi_A(a_2)} \quad (13)$$

Furthermore, since  $u_S$  is strictly supermodular,  $\pi$  is also comonotone by Theorem 1 and Lemma 1, so any  $\theta \in I_1$  and  $\theta' \in I_2$  satisfies  $\theta \leq \theta'$ . Since  $u_R$  is submodular, we have

$u_R(\theta, a_1) - u_R(\theta, a_2) \leq u_R(\theta', a_1) - u_R(\theta', a_2)$  for all  $\theta \in I_1$  and  $\theta' \in I_2$ , which implies

$$\max_{\theta \in I_1} \{u_R(\theta, a_1) - u_R(\theta, a_2)\} \leq \min_{\theta' \in I_2} \{u_R(\theta', a_1) - u_R(\theta', a_2)\}.$$

So

$$\begin{aligned} \sum_{\theta \in I_1} [u_R(\theta, a_1) - u_R(\theta, a_2)] \frac{\pi(\theta, a_1)}{\pi_A(a_1)} &\leq \max_{\theta \in I_1} \{u_R(\theta, a_1) - u_R(\theta, a_2)\} \\ &\leq \min_{\theta' \in I_2} \{u_R(\theta', a_1) - u_R(\theta', a_2)\} \\ &\leq \sum_{\theta' \in I_2} [u_R(\theta', a_1) - u_R(\theta', a_2)] \frac{\pi(\theta', a_2)}{\pi_A(a_2)} \end{aligned} \quad (14)$$

Combining (13) and (14), we have

$$\sum_{\theta \in I_1} [u_R(\theta, a_1) - u_R(\theta, a_2)] \frac{\pi(\theta, a_1)}{\pi_A(a_1)} = \max_{\theta \in I_1} \{u_R(\theta, a_1) - u_R(\theta, a_2)\} = 0$$

and

$$\sum_{\theta' \in I_2} [u_R(\theta', a_1) - u_R(\theta', a_2)] \frac{\pi(\theta', a_2)}{\pi_A(a_2)} = \min_{\theta' \in I_2} \{u_R(\theta', a_1) - u_R(\theta', a_2)\} = 0$$

So  $u_R(\theta, a_1) = u_R(\theta, a_2)$  for all  $\theta \in I_1 \cup I_2$ .

Since the argument above applies to any  $a_1, a_2 \in \text{supp}(\pi_A)$ , we have that for all  $a_i, a, a' \in \text{supp}(\pi_A)$ ,

$$u_R(\theta, a_i) - u_R(\theta, a) = u_R(\theta, a_i) - u_R(\theta, a') = 0 \quad \forall \theta \in I_i,$$

so for all  $i$  and  $\theta \in I_i$ , we have

$$u_R(\theta, a) - u_R(\theta, a') = 0 \quad \forall a, a' \in \text{supp}(\pi_A),$$

and therefore

$$u_R(\theta, a) - u_R(\theta, a') = 0 \quad \forall a, a' \in \text{supp}(\pi_A) \text{ and } \theta \in \Theta.$$

However, this is a contradiction since by assumption, there exists no  $a, a' \in A$  such that  $a \neq a'$  and  $u_R(\theta, a) = u_R(\theta, a')$  for all  $\theta$ .

Therefore  $\text{supp}(\pi_A)$  must be a singleton, denoted by  $a^*$ . Then  $u_R$ -obedience implies  $a^* \in \arg \max_{a \in A} \sum_{\theta} \mu_0(\theta) u(\theta, a)$ . So  $\pi$  is a no-information outcome.

## A.6 Proof of Proposition 2

*Proof of statement 1.* For each  $a \in A$ , let

$$P_a \equiv \{\mu \in \Delta(\Theta) \mid \sum_{\theta} \mu(\theta) u_R(\theta, a) > \sum_{\theta} \mu(\theta) u_R(\theta, a'), \forall a' \neq a\}$$

which denotes the set of beliefs such that  $a$  is the Receiver's strict best response. We prove our claim under the assumption that there exists  $a^\circ \in A$  such that  $\mu_0 \in P_{a^\circ}$  (i.e.  $a^\circ$  is the unique best response to  $\mu_0$ ). Later we will show that this assumption holds for generic priors.

When the Sender's information structure is uninformative, the Receiver best responds to the Sender's messages by choosing  $a^\circ$ . The Sender's payoff is

$$v_0 \equiv \sum_{\theta \in \Theta} \mu_0(\theta) u_S(\theta, a^\circ).$$

We will show that there exists a stable outcome distribution that gives the Sender a higher payoff than  $v_0$ .

We consider the case where the Sender benefits from persuasion, so  $a^\circ \neq \bar{a}$ , otherwise the Receiver is already choosing the Sender's favorite action under the prior. For  $\varepsilon$  sufficiently small, consider the outcome distribution  $\pi^\varepsilon \in \Delta(\Theta \times A)$  defined by

$$\pi^\varepsilon(\theta, a) = \begin{cases} \mu_0(\theta) & \text{if } \theta \neq \bar{\theta}, a = a^\circ \\ \mu_0(\bar{\theta}) - \varepsilon & \text{if } (\theta, a) = (\bar{\theta}, a^\circ) \\ \varepsilon & \text{if } (\theta, a) = (\bar{\theta}, \bar{a}) \\ 0 & \text{otherwise .} \end{cases}$$

We will show that for  $\varepsilon$  sufficiently small,  $\pi^\varepsilon$  is stable and gives the Sender a higher payoff than  $v_0$ .

It can be easily seen that the support of  $\pi^\varepsilon$  is comonotone. Since  $u_S$  is supermodular,  $\pi^\varepsilon$  is  $u_S$ -cyclically monotone by Lemma 1. Next, we verify that for  $\varepsilon$  sufficiently small,  $\pi^\varepsilon$  satisfies  $u_R$ -obedience at the two actions  $\{\bar{a}, a^\circ\}$ . For  $a^\circ$ , note that since  $\mu_0 \in P_{a^\circ}$ , we have

$$\sum_{\theta \in \Theta} \mu_0(\theta) u(\theta, a^\circ) > \sum_{\theta \in \Theta} \mu_0(\theta) \pi(\theta, a') \text{ for all } a' \in A,$$

so for  $\varepsilon$  sufficiently small,

$$\sum_{\theta \in \Theta} \mu_0(\theta) u(\theta, a^\circ) - \varepsilon u(\bar{\theta}, a^\circ) \geq \sum_{\theta \in \Theta} \mu_0(\theta) \pi(\theta, a') - \varepsilon u(\bar{\theta}, a') \text{ for all } a' \in A.$$

which means  $\pi^\varepsilon$  satisfies  $u_R$ -obedience at  $a^\circ$ .

As  $\bar{a} \in A^\circ$ , there exists  $\bar{\mu} \in \Delta(\Theta)$  such that  $\bar{a} \in \arg \max_a \sum_\theta \bar{\mu}(\theta) u_R(\theta, a)$ . So for every  $a' \neq \bar{a}$ ,

$$\sum_\theta \bar{\mu}(\theta) [u_R(\theta, \bar{a}) - u_R(\theta, a')] \geq 0$$

Since  $u_R$  is supermodular,  $u_R(\theta, \bar{a}) - u_R(\theta, a')$  is weakly increasing in  $\theta$ , so if a belief  $\mu'$  first order stochastically dominates  $\bar{\mu}$ , then

$$\sum_\theta \mu'(\theta) [u_R(\theta, \bar{a}) - u_R(\theta, a')] \geq \sum_\theta \bar{\mu}(\theta) [u_R(\theta, \bar{a}) - u_R(\theta, a')] \geq 0 \text{ for all } a' \neq \bar{a}.$$

In particular, the Dirac measure  $\delta_{\bar{\theta}}$  first order stochastically dominates  $\bar{\mu}$ , so the inequality above implies

$$u_R(\bar{\theta}, \bar{a}) - u_R(\bar{\theta}, a') \geq 0 \text{ for all } a' \neq \bar{a}.$$

So  $\bar{a} \in \arg \max_a u_R(\bar{\theta}, a)$ , and  $\pi^\varepsilon$  is  $u_R$ -obedient at action  $\bar{a}$ .

Finally, we show that the Sender obtains higher payoff from  $\pi^\varepsilon$  than  $v_0$ . Note that since by our assumption,  $u_S(\bar{\theta}, a') < u_S(\bar{\theta}, \bar{a})$  for all  $a' \neq \bar{a}$ , we have

$$\begin{aligned} \sum_{\theta, a} \pi^\varepsilon(\theta, a) u_S(\theta, a) &= \sum_{\theta \neq \bar{\theta}} \mu_0(\theta) u_S(\theta, a^\circ) + (\mu_0(\bar{\theta}) - \varepsilon) u_S(\bar{\theta}, a^\circ) + \varepsilon u_S(\bar{\theta}, \bar{a}) \\ &> \sum_{\theta \neq \bar{\theta}} \mu_0(\theta) u_S(\theta, a^\circ) + (\mu_0(\bar{\theta}) - \varepsilon) u_S(\bar{\theta}, a^\circ) + \varepsilon u_S(\bar{\theta}, a^\circ) \\ &= \sum_\theta \mu_0(\theta) u_S(\theta, a^\circ) = v_0. \end{aligned}$$

Therefore, Sender receives a strictly higher payoff from  $\pi^\varepsilon$  than  $v_0$ . This completes the proof.

The rest of the proof shows that  $\cup_{a \in A} P_a$  contains an open, dense, and full measure set. In particular, we show that  $\Delta(\Theta)/\{\cup_{a \in A} P_a\}$  is included in a negligible, closed, and nowhere dense subset in  $\Delta(\Theta)$ .

Define  $H_{a,a'} \equiv \{\mu \in \Delta(\Theta) \mid \sum_\theta \mu(\theta) (u_R(\theta, a) - u_R(\theta, a')) = 0\}$  for any  $a \neq a'$ . Since by [Assumption 1](#),  $u_R(\cdot, a) - u_R(\cdot, a') \neq \mathbf{0}$ , which implies  $J_{a,a'} \equiv \{\mu \in \mathbb{R}^{|\Theta|} \mid \sum_\theta \mu(\theta) (u_R(\theta, a) - u_R(\theta, a')) = 0\}$  is a hyperplane in  $\mathbb{R}^{|\Theta|}$ . Notice that  $H_{a,a'} = J_{a,a'} \cap \Delta(\Theta)$ , which is the intersection of a hyperplane with a simplex. Since the hyperplane includes  $\mathbf{0}$  and  $\Delta(\Theta)$  doesn't, they have to either be parallel with no intersection, or their intersection is in a lower dimensional subspace, which is negligible, closed, and nowhere dense.

For any  $\mu \in \Delta(\Theta)/\{\cup_{a \in A} P_a\}$ , since the maximizer of  $\sum_\theta \mu(\theta) u_R(\theta, a)$  is not unique, there exists  $a, a'$  such that  $\sum_\theta \mu(\theta) (u_R(\theta, a) - u_R(\theta, a')) = 0$ . So  $\Delta(\Theta)/\{\cup_{a \in A} P_a\} \subseteq \cup_{a \neq a'} H_{a,a'}$ , where the latter is a negligible, closed and nowhere dense subset of  $\Delta(\Theta)$ .

□

*Proof of statement 2.* For any generic prior  $\mu^\circ \in \cup_{a \in A} P_a$ , either  $\mu^\circ \notin P_{\underline{a}}$  or  $\mu^\circ \notin P_{\bar{a}}$ . We consider the case  $\mu^\circ \notin P_{\bar{a}}$ , and the other case can be shown symmetrically. Similarly to the previous argument, for  $\varepsilon$  sufficiently small, consider the outcome distribution  $\pi^\varepsilon \in \Delta(\Theta \times A)$ :

$$\pi^\varepsilon(\theta, a) = \begin{cases} \mu_0(\theta) & \text{if } \theta \neq \bar{\theta}, a = a^\circ \\ \mu_0(\bar{\theta}) - \varepsilon & \text{if } (\theta, a) = (\bar{\theta}, a^\circ) \\ \varepsilon & \text{if } (\theta, a) = (\bar{\theta}, \bar{a}) \\ 0 & \text{otherwise} \end{cases}$$

As we have shown in the proof of statement 1, for  $\varepsilon$  sufficiently small,  $\pi^\varepsilon$  is stable, and gives the Sender a higher payoff than  $v_0$ . Therefore, the Sender benefits from credible persuasion. □

*Proof of statement 3.* Let  $\Pi_F$  denote the set of fully revealing outcome distributions, which is compact because it is a closed subset of  $\Delta(\Theta \times A)$ . Let

$$\Pi_F^* \equiv \arg \max_{\pi \in \Pi_F} \sum_{\theta, a} \pi(\theta, a) u_S(\theta, a)$$

be the subset of  $\Pi_F$  that maximizes Sender's payoff, which is also compact by Berge's theorem of maximum. Note that by definition, every fully-revealing outcome distribution is obedient. We will show that there exists an outcome distribution  $\pi^* \in \Pi_F^*$  that is comonotone. This implies that as long as the Sender benefits from one fully-revealing outcome distribution, she must also benefit from  $\pi^*$ , which is a comonotone (and obedient) fully-revealing outcome distribution. This will then complete our proof following [Theorem 1](#) and [Lemma 1](#).

To this end, let us choose

$$\pi^* \in \arg \max_{\pi \in \Pi_F^*} \sum_{\theta, a} \pi(\theta, a) \theta a.$$

Suppose by contradiction that  $\pi^*$  is not comonotone, we construct another outcome distribution  $\pi' \in \Pi_F^*$  that satisfies  $\sum_{\theta, a} \pi'(\theta, a) \theta a > \sum_{\theta, a} \pi^*(\theta, a) \theta a$ , which contradicts  $\pi^* \in \arg \max_{\pi \in \Pi_F^*} \sum_{\theta, a} \pi(\theta, a) \theta a$ .

Since  $\pi^*$  is not comonotone, there exists a pair  $(\theta_1, a_1), (\theta_2, a_2)$  in the support of  $\pi^*$  such that  $\theta_1 < \theta_2$  and  $a_1 > a_2$ . Take  $\varepsilon = \min\{\pi^*(\theta_1, a_1), \pi^*(\theta_2, a_2)\}$ , and construct the outcome distribution  $\pi'$  where:

- $\pi'(\theta_1, a_1) = \pi^*(\theta_1, a_1) - \varepsilon$ ,  $\pi'(\theta_2, a_2) = \pi^*(\theta_2, a_2) - \varepsilon$



- $\pi'(\theta_1, a_2) = \pi^*(\theta_1, a_2) + \varepsilon$ ,  $\pi'(\theta_2, a_1) = \pi^*(\theta_2, a_1) + \varepsilon$
- $\pi'(\theta, a) = \pi^*(\theta, a)$  for all other  $(\theta, a)$

We first argue that  $\pi' \in \Pi_F^*$ . Let  $A^*(\theta) \equiv \arg \max_{a \in A} u_R(\theta, a)$  denote the Receiver's best response correspondence. Since  $u_R(\theta, a)$  is supermodular, by Lemma 2.8.1 of [Topkis \(2011\)](#),  $A^*(\theta)$  is weakly increasing in  $\theta$  in the induced set order. That is, for any  $\theta > \theta'$ ,  $a \in A^*(\theta)$ , and  $a' \in A^*(\theta')$ , we have  $\max\{a, a'\} \in A^*(\theta)$  and  $\min\{a, a'\} \in A^*(\theta')$ . Since  $a_1 \in A^*(\theta_1)$  and  $a_2 \in A^*(\theta_2)$ , we have  $a_1 \in A^*(\theta_2)$  and  $a_2 \in A^*(\theta_1)$ . Therefore,  $\pi'$  is also a fully revealing outcome distribution. Moreover since  $u_S$  is supermodular,

$$\sum_{\theta, a} [\pi'(\theta, a) - \pi^*(\theta, a)] u_S(\theta, a) = \varepsilon [u_S(\theta_1, a_2) + u_S(\theta_2, a_1) - u_S(\theta_1, a_1) - u_S(\theta_2, a_2)] \geq 0,$$

so the Sender's payoff from  $\pi'$  is weakly greater than from  $\pi^*$ , and therefore  $\pi' \in \Pi_F^*$ .

Next we argue that  $\sum_{\theta, a} \pi'(\theta, a) \theta a > \sum_{\theta, a} \pi^*(\theta, a) \theta a$ . To this end, note that

$$\begin{aligned} \sum_{\theta, a} [\pi'(\theta, a) - \pi^*(\theta, a)] \theta a &= \varepsilon [\theta_1 a_2 + \theta_2 a_1 - \theta_1 a_1 - \theta_2 a_2] \\ &= (\theta_2 - \theta_1)(a_1 - a_2) > 0. \end{aligned}$$

This contradicts  $\pi^* \in \arg \max_{\pi \in \Pi_F^*} \sum_{\theta, a} \pi(\theta, a) \theta a$ . □

## A.7 Proof of [Proposition 3](#)

From Theorem 1 of [Mensch \(2021\)](#), if both  $u_S$  and  $u_R$  are supermodular and  $|A| = 2$ , there exists a KG optimal outcome distribution that is comonotone. Then by [Theorem 1](#) and [Lemma 1](#), such an outcome distribution is stable. Moreover, if in addition  $u_S$  is strictly supermodular, any KG optimal outcome distribution is comonotone. So any KG optimal outcome distribution is stable.

## A.8 Proof of [Proposition 4](#)

We begin by establishing a lemma that will be useful for proving [Proposition 4](#).

**Lemma 4.** *Suppose the message space  $M$  is a finite subset of  $\mathbb{R}$ , the information structure  $\lambda \in \Delta(\Theta \times M)$  is comonotone, and the Receiver's payoff function  $u_R$  is strictly supermodular. Consider a Receiver strategy  $\sigma : M \rightarrow A$  defined by*

$$\sigma(m) \in \arg \max_{a \in A} \sum_{\theta} \lambda(\theta, m) u_R(\theta, a).$$

The outcome distribution  $\pi \in \Delta(\Theta \times A)$  induced by  $(\lambda, \sigma)$  is comonotone and  $u_R$ -obedient.

*Proof.* The fact that  $\pi$  is  $u_R$ -obedient follows from [Bergemann and Morris \(2016\)](#). We will prove that  $\pi$  is comonotone. Suppose by contradiction that  $\pi$  is not comonotone, so there exists  $(\theta_1, a_1), (\theta_2, a_2) \in \text{supp}(\pi)$ , such that  $a_1 > a_2$  and  $\theta_1 < \theta_2$ . We will show that this leads to a contradiction.

Let  $M_1 = \{m \in M : \lambda(\theta_1, m) > 0\}$  and  $M_2 = \{m \in M : \lambda(\theta_2, m) > 0\}$ . Since  $(\theta_1, a_1), (\theta_2, a_2) \in \text{supp}(\pi)$ , there exists  $m_1 \in M_1$  and  $m_2 \in M_2$  such that  $\sigma(m_1) = a_1$  and  $\sigma(m_2) = a_2$ . In addition,  $m_1 \leq m_2$  because  $\theta_1 < \theta_2$  and  $\lambda$  is comonotone; furthermore,  $m_1 \neq m_2$  because  $\sigma(m_1) \neq \sigma(m_2)$ , so  $m_1 < m_2$ .

Let  $\Theta_1 = \{\theta \in \Theta : \lambda(\theta, m_1) > 0\}$  and  $\Theta_2 = \{\theta \in \Theta : \lambda(\theta, m_2) > 0\}$ . Since  $\sigma$  best responds to each message, we have

$$\sum_{\theta \in \Theta_1} [u_R(\theta, a_1) - u_R(\theta, a_2)] \frac{\lambda(\theta, m_1)}{\lambda_M(m_1)} \geq 0 \geq \sum_{\theta' \in \Theta_2} [u_R(\theta', a_1) - u_R(\theta', a_2)] \frac{\lambda(\theta', m_2)}{\lambda_M(m_2)} \quad (15)$$

Furthermore, since  $\lambda$  is comonotone and  $m_1 < m_2$ , for any  $\theta \in \Theta_1$  and  $\theta' \in \Theta_2$ ,  $\theta \leq \theta'$ , which implies  $\max \Theta_1 \leq \min \Theta_2$ . Together with the supermodularity of  $u_R$ , we have

$$\max_{\theta \in \Theta_1} \{u_R(\theta, a_1) - u_R(\theta, a_2)\} \leq \min_{\theta' \in \Theta_2} \{u_R(\theta', a_1) - u_R(\theta', a_2)\}.$$

So

$$\begin{aligned} \sum_{\theta \in \Theta_1} [u_R(\theta, a_1) - u_R(\theta, a_2)] \frac{\lambda(\theta, m_1)}{\lambda_M(m_1)} &\leq \max_{\theta \in \Theta_1} \{u_R(\theta, a_1) - u_R(\theta, a_2)\} \\ &\leq \min_{\theta' \in \Theta_2} \{u_R(\theta', a_1) - u_R(\theta', a_2)\} \\ &\leq \sum_{\theta' \in \Theta_2} [u_R(\theta', a_1) - u_R(\theta', a_2)] \frac{\lambda(\theta', m_2)}{\lambda_M(m_2)} \end{aligned} \quad (16)$$

Combining (15) and (16), we have

$$\sum_{\theta \in \Theta_1} [u_R(\theta, a_1) - u_R(\theta, a_2)] \frac{\lambda(\theta, m_1)}{\lambda_M(m_1)} = \max_{\theta \in \Theta_1} \{u_R(\theta, a_1) - u_R(\theta, a_2)\} = 0$$

and

$$\sum_{\theta' \in \Theta_2} [u_R(\theta', a_1) - u_R(\theta', a_2)] \frac{\lambda(\theta', m_2)}{\lambda_M(m_2)} = \min_{\theta' \in \Theta_2} \{u_R(\theta', a_1) - u_R(\theta', a_2)\} = 0$$

so  $u_R(\theta, a_1) = u_R(\theta, a_2)$  for all  $\theta \in \Theta_1 \cup \Theta_2$ .

But recall that  $\theta_1, \theta_2 \in \Theta_1 \cup \Theta_2$  and  $\theta_1 < \theta_2$ , and from the strict supermodularity of  $u_R$ ,

$$u_R(\theta_1, a_1) - u_R(\theta_1, a_2) < u_R(\theta_2, a_1) - u_R(\theta_2, a_2),$$

which leads to a contradiction.  $\square$

*Proof of Proposition 4.* Let  $\pi$  be a Sender-optimal stable outcome distribution under preferences  $(u_S, u_R)$ . By Theorem 1 and Lemma 1,  $\pi$  is comonotone.

Now under the more aligned preferences  $(u_S, u'_R)$ , suppose the Sender uses the information structure  $\lambda = \pi$  with message space  $M = \text{supp}(\pi_A)$ , and let  $\sigma'$  be the Receiver strategy that best responds to each message from  $\pi$ , with ties broken in favor of the Sender. By Lemma 4, the outcome distribution  $\pi'$  induced by the profile  $(\pi, \sigma')$  is comonotone and  $u'_R$ -obedient. By Theorem 1 and Lemma 1,  $\pi'$  is a stable outcome distribution under preferences  $(u_S, u'_R)$ .

It remains to show that the Sender obtains a higher payoff from  $\pi'$ . For each belief  $\mu \in \Delta(\Theta)$ ,

$$\hat{a}(\mu) \in \arg \max_{a \in A} \sum_{\theta} \mu(\theta) u_R(\theta, a) \quad \text{and} \quad \hat{a}'(\mu) \in \arg \max_{a \in A} \sum_{\theta} \mu(\theta) u'_R(\theta, a)$$

denote the Receiver's best response to belief  $\mu$ , with ties broken in favor of the Sender. Note that since  $\sigma'$  breaks ties in favor of the Sender,

$$E_{\pi(\cdot|a)}[u_S(\theta, \sigma(a))] = E_{\pi(\cdot|a)}[u_S(\theta, \hat{a}'(\pi(\cdot|a)))] \quad \text{for all } a \in M. \quad (17)$$

By contrast,

$$E_{\pi(\cdot|a)}[u_S(\theta, a)] \leq E_{\pi(\cdot|a)}[u_S(\theta, \hat{a}(\pi(\cdot|a)))] \quad \text{for all } a \in M \quad (18)$$

since  $\pi$  may not be the result of a Sender-favoring tie-breaking strategy.

So

$$\begin{aligned} E_{\pi'}[u_S(\theta, a)] &= E_{\pi}[u_S(\theta, \sigma(a))] \\ &= E_{\pi_A} \left[ E_{\pi(\cdot|a)}[u_S(\theta, \sigma(a))] \right] \\ &= E_{\pi_A} \left[ E_{\pi(\cdot|a)}[u_S(\theta, \hat{a}'(\pi(\cdot|a)))] \right] \\ &\geq E_{\pi_A} \left[ E_{\pi(\cdot|a)}[u_S(\theta, \hat{a}(\pi(\cdot|a)))] \right] \\ &\geq E_{\pi_A} \left[ E_{\pi(\cdot|a)}[u_S(\theta, a)] \right] \\ &= E_{\pi}[u_S(\theta, a)] \end{aligned}$$

where the first line above follows from the definition of  $\pi'$ , the second line is the law of iterated expectation, the third follows from (17), the fourth line follows from the preferences  $(u_S, u'_R)$  being more aligned than  $(u_S, u_R)$ , the fifth follows from (18), and the last equality is again the law of iterated expectation.

So the Sender obtains a higher payoff from  $\pi'$  than  $\pi$ . Since  $\pi'$  is also stable under preferences  $(u_S, u'_R)$ , this completes our proof.  $\square$

## A.9 Proof of Proposition 5

For each buyer's belief over quality,  $\mu \in \Delta(\Theta)$ , let  $\underline{\theta}_\mu$  denote the smallest  $\theta$  in the support of  $\mu$ ; in addition, let  $\phi_\mu(x) \equiv E_\mu[v(\theta)|\theta \leq x]$  denote the corresponding expected value to buyers when the quality threshold is  $\theta \leq x$ .<sup>27</sup> Clearly,  $\phi_\mu(\cdot)$  is weakly increasing and  $\phi_\mu(1) = E_\mu[v(\theta)]$ .

**Lemma 5.** *For every  $\mu \in \Delta(\Theta)$ , there exists a largest fixed point  $\theta_\mu^* \in (\underline{\theta}_\mu, 1)$  such that  $\phi_\mu(\theta_\mu^*) = \theta_\mu^*$ . Moreover, for any  $\theta \in (\theta_\mu^*, 1]$ ,  $\phi_\mu(\theta) < \theta$ .*

*Proof.* Since  $\phi_\mu(\underline{\theta}_\mu) = v(\underline{\theta}_\mu) > \underline{\theta}_\mu$ ,  $\phi_\mu(1) = E_\mu[v(\theta)] < 1$ , and  $\phi_\mu(\cdot)$  is weakly increasing, from Tarski's fixed point theorem, there exists a largest fixed point  $\theta_\mu^* \in (\underline{\theta}_\mu, 1)$  such that  $\phi_\mu(\theta_\mu^*) = \theta_\mu^*$ . To see the second statement, suppose there exists  $\theta \in (\theta_\mu^*, 1)$  such that  $\phi_\mu(\theta) \geq \theta$ , again from Tarski's fixed point theorem, there exists a fixed point  $\theta' \in (\theta_\mu^*, 1)$ , which contradicts  $\theta_\mu^*$  being the largest fixed point.  $\square$

**Lemma 6.** *Let  $\lambda \in \Delta(\Theta \times M)$  be an information structure, and for every  $m \in M$  let  $\mu_m \in \Delta(\Theta)$  denote the buyers' posterior belief after observing message  $m$ . The following strategy profile is a BNE in the game  $\langle G, \lambda \rangle$ :  $\alpha_S(\theta, m) = \theta$ ,  $\beta_1(m) = \beta_2(m) = \theta_{\mu_m}^*$ .*

*Proof.* For every message  $m$ , since  $\phi_{\mu_m}(\theta_{\mu_m}^*) = \theta_{\mu_m}^*$ , each buyer's expected payoff is 0. Any deviation to a lower bid also gives a payoff of zero. From Lemma 5, for any  $\theta \in (\theta_{\mu_m}^*, 1]$ ,  $\phi_{\mu_m}(\theta) < \theta$ , so any deviation to a bid higher than  $\theta_{\mu_m}^*$  would lead to a negative payoff. Therefore no buyer has an incentive to deviate.  $\square$

**Lemma 7.** *Let  $(\lambda^*, \sigma^*)$  be a WD-IC profile. For each message  $m$ , let  $p(m) \equiv \max\{\beta_1^*(m), \beta_2^*(m)\}$  denote the equilibrium market price in the game  $\langle G, \lambda^* \rangle$ . Then  $\phi_{\mu_m}(p(m)) = p(m)$  for each  $m \in M$ .*

*Proof.* Suppose  $\phi_{\mu_m}(p(m)) < p(m)$ , then the winning buyer's payoff is negative, and can profitably deviate to bid 0. Now suppose  $\phi_{\mu_m}(p(m)) > p(m)$ , we show that at least one buyer has an incentive to bid a higher price.

If  $\beta_1^*(m) \neq \beta_2^*(m)$ , then the losing bidder can profitably deviate. Since  $\phi_{\mu_m}(\cdot)$  is weakly increasing, there exists small enough  $\varepsilon$  such that  $\phi_{\mu_m}(p(m) + \varepsilon) > p(m) + \varepsilon$ . So the losing bidder can deviate to bidding  $p(m) + \varepsilon$  and receives a strictly positive payoff.

---

<sup>27</sup>For  $x$  less than  $\underline{\theta}_\mu$  we set  $\phi_\mu(x) = v(\underline{\theta}_\mu)$ .

If  $\beta_1^*(m) = \beta_2^*(m) = b$  for some  $b$ , we show that both buyers have an incentive to deviate. Let  $K \equiv \phi_{\mu_m}(b) - b > 0$ . Since ties are broken evenly, each buyer's payoff is  $\frac{1}{2}P_{\mu_m}(\theta \leq b)K$ . By letting  $\varepsilon < \frac{K}{2}$ , we have

$$\phi_{\mu_m}(b + \varepsilon) - b - \varepsilon \geq \phi_{\mu_m}(b) - b - \varepsilon = K - \varepsilon > \frac{K}{2}.$$

So if either of the bidders deviates to bidding  $b + \varepsilon$ , he receives a payoff of  $P_{\mu_m}(\theta \leq b + \varepsilon)[\phi_{\mu_m}(b + \varepsilon) - b - \varepsilon] > \frac{1}{2}P_{\mu_m}(\theta \leq b)K$ , which is profitable.  $\square$

**Lemma 8.** *If a profile  $(\lambda^*, \sigma^*)$  is credible and WD-IC, then there exists a set  $E \subset \Theta \times M$  such that  $\lambda^*(E) = 1$ , and for any  $(\theta, m), (\theta', m') \in E$ ,*

$$\max\{\theta, p(m)\} + \max\{\theta', p(m')\} \geq \max\{\theta, p(m')\} + \max\{\theta', p(m)\}.$$

*Proof.* Since  $(\lambda^*, \sigma^*)$  is WD-IC, trade only happens when the seller's ask price  $\alpha^*(\theta, m) = \theta$  is higher than the prevailing market price  $p(m) = \max\{\beta_1^*(m), \beta_2^*(m)\}$ . The seller's payoff function can therefore be simplified as

$$u_S(\theta, \sigma^*(\theta, m)) = u_S(\theta, \alpha^*(\theta, m), \beta_1^*(m), \beta_2^*(m)) = \max\{\theta, p(m)\}.$$

Recall that credibility requires

$$\lambda \in \arg \max_{\lambda' \in D(\lambda)} \int u_S(\theta, \sigma^*(\theta, m)) d\lambda'(\theta, m).$$

Let  $v_S(\theta, m) \equiv u_S(\theta, \sigma^*(\theta, m)) = \max\{\theta, p(m)\}$ . From Theorem 1 of [Beiglböck et al. \(2009\)](#),  $\lambda$  is  $v_S$ -cyclically monotone. That is, there exists a set  $E \subset \Theta \times M$  such that  $\lambda^*(E) = 1$ , and for any sequence  $(\theta_k, m_k)_{k=1}^n \in E$ ,

$$\sum_{k=1}^n v_S(\theta_k, m_k) \geq \sum_{k=1}^n v_S(\theta_k, m_{k+1}).$$

Suppose  $(\theta, m), (\theta', m') \in E$ , then  $v_S$ -cyclical monotonicity implies that

$$v_S(\theta, m) + v_S(\theta', m') \geq v_S(\theta, m') + v_S(\theta', m),$$

which is

$$\max\{\theta, p(m)\} + \max\{\theta', p(m')\} \geq \max\{\theta, p(m')\} + \max\{\theta', p(m)\}.$$

$\square$

In light of [Lemma 8](#), for every credible profile  $(\lambda^*, \sigma^*)$  we will focus only on pairs  $(\theta, m) \in E$ . We will use  $\text{proj}_M(E) \equiv \{m \in M : (\theta, m) \in E\}$  to denote the projection of  $E$  onto the message space.

Let  $\underline{p} = \inf\{p(m) | m \in \text{proj}_M(E)\}$  be the infimum of trading prices across all messages. For each message  $m$ , let  $\Theta(m) = \{\theta : (\theta, m) \in E\}$  be the set of  $\theta$  that is matched with  $m$ .

**Lemma 9.** *Let  $(\lambda^*, \sigma^*)$  be a credible and WD-IC profile. For every message  $\hat{m} \in \text{proj}_M(E)$  such that  $\hat{p} \equiv p(\hat{m}) = \max\{\beta_1^*(\hat{m}), \beta_2^*(\hat{m})\} > \underline{p}$ , we have  $\Theta(\hat{m}) \cap (\underline{p}, \infty) = \emptyset$ .*

*Proof.* To prove the lemma, suppose by contradiction that there exists  $\hat{\theta} \in \Theta(\hat{m}) \cap (\underline{p}, \infty)$ . By the definition of  $\underline{p}$ , there exists  $p'$  with  $\underline{p} < p' < \hat{\theta}$  such that  $p' = p(m')$  for some  $m' \in \text{proj}_M(E)$ . Since in equilibrium  $p' = E_{\mu_{m'}}[v(\theta) | \theta \in \Theta(m') \cap [0, p']]$ , there also exists  $\theta' \in \Theta(m')$  such that  $\theta' < p'$ . Since  $(\theta', m'), (\hat{\theta}, \hat{m}) \in E$ , by [Lemma 8](#), we have

$$\max\{\theta', \hat{p}\} + \max\{\hat{\theta}, p'\} \leq \max\{\theta', p'\} + \max\{\hat{\theta}, \hat{p}\}$$

Since  $\theta' < p'$  by construction, we have

$$\max\{\theta', \hat{p}\} + \max\{\hat{\theta}, p'\} \leq p' + \max\{\hat{\theta}, \hat{p}\} \quad (19)$$

Note also that

$$p' + \max\{\hat{\theta}, \hat{p}\} < \hat{p} + \hat{\theta}. \quad (20)$$

The inequality above follows by considering two possibilities for  $\max\{\hat{\theta}, \hat{p}\}$ : either  $\hat{\theta} \geq \hat{p}$ , in which case  $p' + \max\{\hat{\theta}, \hat{p}\} = p' + \hat{\theta} < \hat{p} + \hat{\theta}$ ; or  $\hat{\theta} < \hat{p}$ , in which case  $p' + \max\{\hat{\theta}, \hat{p}\} = p' + \hat{p} < \hat{p} + \hat{\theta}$  as well.

Combining (19) and (20), and noticing  $\theta' < p' < \hat{p}$  and  $p' < \hat{\theta}$  yield

$$\max\{\theta', \hat{p}\} + \max\{\hat{\theta}, p'\} < \hat{p} + \hat{\theta} = \max\{\theta', \hat{p}\} + \max\{\hat{\theta}, p'\}$$

which is a contradiction. □

*Proof of the [Proposition 5](#).* To prove our result, we first calculate the seller's profit from an arbitrary credible and WD-IC profile  $(\lambda^*, \sigma^*)$ . We then show that there exists another credible and WD-IC profile  $(\lambda^0, \sigma^0)$ , where  $\lambda^0$  is a null information structure, that leads to weakly higher profit for the seller.

Recall that seller's payoff function can be written as  $u_S(\theta, \sigma^*(\theta, m)) = \max\{\theta, p(m)\}$ , so

her ex-ante profit is

$$\begin{aligned}
& \int_{\Theta \times M} \max\{\theta, p(m)\} d\lambda^*(\theta, m) = \int_M \int_0^1 \max\{\theta, p(m)\} d\lambda^*(\theta|m) d\lambda_M^*(m) \\
&= \int_M \left[ \int_0^{p(m)} p(m) d\lambda^*(\theta|m) + \int_{p(m)}^1 \theta d\lambda^*(\theta|m) \right] d\lambda_M^*(m) \\
&= \int_M \left[ p(m) P_{\lambda^*(\theta|m)}(\theta \leq p(m)) + \int_{p(m)}^1 \theta d\lambda^*(\theta|m) \right] d\lambda_M^*(m)
\end{aligned}$$

By Lemma 7,  $p(m) = E_{\lambda^*(\theta|m)}[v(\theta)|\theta \leq p(m)]$ , so we can write the integral above as

$$\begin{aligned}
& \int_M \left[ E_{\lambda^*(\theta|m)}[v(\theta)|\theta \leq p(m)] P_{\lambda^*(\theta|m)}(\theta \leq p(m)) + \int_{p(m)}^1 \theta d\lambda^*(\theta|m) \right] d\lambda_M^*(m) \\
&= \int_M \left[ \int_0^{p(m)} v(\theta) d\lambda^*(\theta|m) + \int_{p(m)}^1 \theta d\lambda^*(\theta|m) \right] d\lambda_M^*(m)
\end{aligned}$$

By Lemma 9, for every  $m \in \text{proj}_M(E)$ , if  $p(m) > \underline{p}$  then  $\Theta(m) \cap (\underline{p}, \infty) = \emptyset$ , so the seller's profit from  $(\lambda^*, \sigma^*)$  can be further simplified to

$$\begin{aligned}
& \int_M \left[ \int_0^{\underline{p}} v(\theta) d\lambda^*(\theta|m) + \int_{\underline{p}}^1 \theta d\lambda^*(\theta|m) \right] d\lambda_M^*(m) \\
&= \int_0^{\underline{p}} v(\theta) d\mu_0(\theta) + \int_{\underline{p}}^1 \theta d\mu_0(\theta).
\end{aligned} \tag{21}$$

Having calculated the seller's profit from  $(\lambda^*, \sigma^*)$ , next we will construct another credible and WD-IC profile  $(\lambda^0, \sigma^0)$  with a weakly higher profit, where  $\lambda_0$  is the null information structure  $\mu_0 \times \delta_{m_0}$ .

From Lemma 9, for every  $m \in \text{proj}_M(E)$ ,

$$\phi_{\mu_m}(p(m)) = E_{\mu_m}[v(\theta)|\theta \leq p(m)] = E_{\mu_m}[v(\theta)|\theta \leq \underline{p}] = \phi_{\mu_m}(\underline{p});$$

in addition, from Lemma 7,  $\phi_{\mu_m}(p(m)) = p(m)$  for every message  $m \in \text{proj}_M(E)$ . Combining these yields

$$\phi_{\mu_m}(\underline{p}) = \phi_{\mu_m}(p(m)) = p(m) \geq \underline{p}.$$

Taking expectation over all messages, we have  $\phi_{\mu_0}(\underline{p}) \geq \underline{p}$ . By Tarski's fixed point theorem, there exists a largest  $p^0 \in [\underline{p}, 1)$  such that the  $\phi_{\mu_0}(p^0) = p^0$ .

Using a similar argument as that in Lemma 6, the strategy profile  $\sigma^0$  where the seller plays her weakly dominant strategy  $\alpha^0(\theta, m_0) = \theta$ , and buyers play  $\beta_1^0(m_0) = \beta_2^0(m_0) = p^0$  is a BNE

in the game  $\langle G, \lambda^0 \rangle$ .

It remains to show that the seller's profit from  $(\lambda^0, \sigma^0)$  is weakly higher than that from  $(\lambda^*, \sigma^*)$ . Under  $(\lambda^0, \sigma^0)$  the seller's profit is

$$\begin{aligned}
\int_0^1 \max\{\theta, p^0\} d\mu_0(\theta) &= \int_0^{p^0} p^0 d\mu_0(\theta) + \int_{p^0}^1 \theta d\mu_0(\theta) \\
&= p^0 P_{\mu_0}(\theta \leq p^0) + \int_{p^0}^1 \theta d\mu_0(\theta) \\
&= E_{\mu_0}[v(\theta) | \theta \leq p^0] P_{\mu_0}(\theta \leq p^0) + \int_{p^0}^1 \theta d\mu_0(\theta) \\
&= \int_0^{p^0} v(\theta) d\mu_0(\theta) + \int_{p^0}^1 \theta d\mu_0(\theta)
\end{aligned} \tag{22}$$

Comparing (21) and (22), since  $p^0 \geq \underline{p}$  and  $v(\theta) > \theta$  for all  $\theta$ , it follows that

$$\int_0^{p^0} v(\theta) d\mu_0 + \int_{p^0}^1 \theta d\mu_0 \geq \int_0^{\underline{p}} v(\theta) d\mu_0 + \int_{\underline{p}}^1 \theta d\mu_0.$$

The seller's profit under  $(\lambda^0, \sigma^0)$  is therefore weakly higher than that from  $(\lambda^*, \sigma^*)$ . □

## A.10 Proof of Proposition 6

For  $\mu \in \Delta(\Theta)$  and  $\nu \in \Delta(A)$ , let  $\Lambda(\mu, \nu) \equiv \{\lambda \in \Delta(\Theta \times A) : \lambda_\Theta = \mu, \lambda_A = \nu\}$  denote the set of joint distributions on  $\Theta \times A$  that with marginals given by  $\mu$  and  $\nu$ . The following lemmas will be useful in our proofs.

**Lemma 10.** *The correspondence*

$$B(\mu, \nu) \equiv \arg \max_{\lambda \in \Lambda(\mu, \nu)} \sum_{\theta, m} \lambda(\theta, m) u_S(\theta, \sigma(m))$$

*is upper hemi-continuous with respect to  $(\mu, \nu)$ . Thus, the value function*

$$V(\mu, \nu) \equiv \max_{\lambda \in \Lambda(\mu, \nu)} \sum_{\theta, m} \lambda(\theta, m) u_S(\theta, \sigma(m))$$

*is continuous.*

*Proof.* The first statement follows directly from Theorem 1.50 of [Santambrogio \(2015\)](#). For any sequence  $(\lambda_k, \mu_k, \nu_k) \rightarrow (\lambda, \mu, \nu)$  so that  $\lambda_k \in B(\mu_k, \nu_k)$  for all  $k$ , we have  $\lambda \in B(\mu, \nu)$ . Then



$V(\mu, \nu) = \sum_{\theta, m} \lambda(\theta, m) u_S(\theta, \sigma(m)) = \lim_{k \rightarrow \infty} \sum_{\theta, m} \lambda_k(\theta, m) u_S(\theta, \sigma(m)) = \lim_{k \rightarrow \infty} V(\nu_k, \nu_k)$ , which proves the second statement.  $\square$

**Lemma 11.** *Suppose  $\mu \in \mathcal{F}_\Theta^N$  and  $\nu \in \mathcal{F}_M^N$ , then the extreme points of  $\Lambda(\mu, \nu)$  are contained in  $X^N(\mu, \nu)$ .*

*Proof.* Consider the set  $Y^N(\mu, \nu) = \{f \in \mathbb{R}_+^{|\Theta| \times |M|} : \sum_{\theta} f(\theta, \cdot) = N\nu(\cdot), \sum_m f(\cdot, m) = N\mu(\cdot)\}$ . From Corollary 8.1.4 of [Brualdi \(2006\)](#), the extreme points of  $Y^N(\mu, \nu)$  are contained in  $Z^N(\mu, \nu) = \{f \in \mathbb{N}^{|\Theta| \times |M|} : \sum_{\theta} f(\theta, \cdot) = N\nu(\cdot), \sum_m f(\cdot, m) = N\mu(\cdot)\}$ . Since  $\Lambda(\mu, \nu) = \{\frac{f}{N} : f \in Y^N(\mu, \nu)\}$  and  $X^N(\mu, \nu) = \{\frac{f}{N} : f \in Z^N(\mu, \nu)\}$ , the extreme points of  $\Lambda(\mu, \nu)$  are contained in  $X^N(\mu, \nu)$ .  $\square$

**Lemma 12.** *Let  $X, Y$  be metric spaces and  $\Gamma : X \rightrightarrows Y$  be a correspondence. If  $\Gamma$  is upper hemicontinuous at  $x_0 \in X$ , and  $\Gamma(x_0) = \{y_0\}$  for some  $y_0 \in Y$ , then  $\Gamma$  is continuous at  $x_0$ .*

*Proof.* For any  $\varepsilon > 0$ , let  $B(y_0, \varepsilon) \subseteq Y$  denote the  $\varepsilon$ -ball centered at  $y_0$ . We will show that there exists  $\delta > 0$  such that for all  $|x - x_0| < \delta$ ,  $\Gamma(x) \cap B(y_0, \varepsilon) \neq \emptyset$ , which implies that  $\Gamma$  is lhc and therefore continuous.

Now since  $\Gamma(x) = \{y_0\} \subseteq B(y_0, \varepsilon)$  and  $\Gamma$  is uhc at  $x_0$ , it follows that there exists  $\delta > 0$  such that  $\Gamma(x) \subseteq B(y_0, \varepsilon)$  for all  $|x - x_0| < \delta$ , so  $\Gamma(x) \cap B(y_0, \varepsilon) \neq \emptyset$  for all  $|x - x_0| < \delta$ , which completes the proof.  $\square$

*Proof of Proposition 6 statement 1.* First suppose  $(\lambda^*, \sigma^*)$  is not credible. Then there exists  $\lambda' \in \Lambda(\mu_0, \lambda_M^*)$  (recall  $\mu_0$  is the prior distribution on  $\Theta$ ) and  $\varepsilon_0 > 0$  such that

$$\sum_{\theta, m} \lambda^*(\theta, m) u_S(\theta, \sigma^*(m)) < \sum_{\theta, m} \lambda'(\theta, m) u_S(\theta, \sigma^*(m)) - \varepsilon_0$$

By continuity, there exists  $\varepsilon_1 > 0$  such that for all  $|\lambda - \lambda^*| < \varepsilon_1$  we have

$$\sum_{\theta, m} \lambda(\theta, m) u_S(\theta, \sigma^*(m)) < \sum_{\theta, m} \lambda'(\theta, m) u_S(\theta, \sigma^*(m)) - \varepsilon_0/2 \quad (23)$$

By [Lemma 10](#), there exists  $\varepsilon_2 > 0$  such that for all  $|\mu - \mu_0| < \varepsilon_2$  and  $|\nu - \lambda_M^*| < \varepsilon_2$ , there exists  $\lambda \in \Lambda(\mu, \nu)$  with

$$\sum_{\theta, m} \lambda(\theta, m) u_S(\theta, \sigma^*(m)) > \sum_{\theta, m} \lambda'(\theta, m) u_S(\theta, \sigma^*(m)) - \varepsilon_0/2. \quad (24)$$

Moreover, since the Receiver is choosing only pure strategies, there exists  $\varepsilon_3$  such that for any  $\sigma$  where  $|\sigma - \sigma^*| < \varepsilon_3$ ,  $\sigma = \sigma^*$ .

Now let  $\varepsilon = \min\{\varepsilon_1, \frac{\varepsilon_2}{|\Theta| \times |M|}, \varepsilon_3\}$ . By assumption, there exists a finite-sample, credible and R-IC profile  $(\nu^N, \phi^N, \sigma^N)$  such that  $|\nu^N - \lambda_M^*| < \varepsilon$ ,  $|\sigma^N - \sigma| < \varepsilon$  and  $P(|\phi^N(F_\Theta^N) - \lambda| < \varepsilon) > 1 - \varepsilon$ .

Under such a finite-sample profile,  $\sigma^N = \sigma^*$ , and there exists  $F_\Theta^N \in \mathcal{F}_\Theta^N$ , realized with positive probability, such that  $\tilde{\lambda}^* = \phi^N(F_\Theta^N)$  satisfies  $|\tilde{\lambda}^* - \lambda^*| < \min\{\varepsilon_1, \frac{\varepsilon_2}{|\Theta| \times |M|}\}$ .

Now since  $|\tilde{\lambda}^* - \lambda^*| < \varepsilon_1$ , by (23) we know that

$$\sum_{\theta, m} \tilde{\lambda}^*(\theta, m) u_S(\theta, \sigma^*(m)) < \sum_{\theta, m} \lambda'(\theta, m) u_S(\theta, \sigma^*(m)) - \varepsilon_0/2 \quad (25)$$

In addition, since  $|\tilde{\lambda}^* - \lambda^*| < \frac{\varepsilon_2}{|\Theta| \times |M|}$ , we know  $F_\Theta^N = \tilde{\lambda}_\Theta^*$  satisfies  $|F_\Theta^N - \mu_0| < \varepsilon_2$ , and  $\nu^N = \tilde{\lambda}_M^*$  satisfies  $|\nu^N - \lambda_M^*| < \varepsilon_2$ , so by (24) there exists  $\tilde{\lambda}' \in \Lambda(F_\Theta^N, \nu^N)$  such that

$$\sum_{\theta, m} \tilde{\lambda}'(\theta, m) u_S(\theta, \sigma^*(m)) > \sum_{\theta, m} \lambda'(\theta, m) u_S(\theta, \sigma^*(m)) - \varepsilon_0/2 \quad (26)$$

Combining (25) and (26), we have

$$\sum_{\theta, m} \tilde{\lambda}'(\theta, m) u_S(\theta, \sigma^*(m)) > \sum_{\theta, m} \tilde{\lambda}^*(\theta, m) u_S(\theta, \sigma^*(m)),$$

Note that  $\tilde{\lambda}' \in \Lambda(F_\Theta^N, \nu^N)$ , but by Lemma 11, we can replace  $\tilde{\lambda}'$  with an extreme point in  $X^N(F_\Theta^N, \nu^N)$ , and the above inequality still holds. That is, there exists  $\hat{\lambda}' \in X^N(F_\Theta^N, \nu^N)$  such that

$$\sum_{\theta, m} \hat{\lambda}'(\theta, m) u_S(\theta, \sigma^*(m)) > \sum_{\theta, m} \tilde{\lambda}^*(\theta, m) u_S(\theta, \sigma^*(m)),$$

Notice that  $\tilde{\lambda}^*$  and  $\hat{\lambda}'$  are both in  $X^N(F_\Theta^N, \nu^N)$ , which is a contradiction since by the credibility of  $(\nu^N, \phi^N, \sigma^N)$

$$\tilde{\lambda}^* = \phi^N(F_\Theta^N) = \arg \max_{\lambda \in X^N(F_\Theta^N, \nu^N)} \sum_{\theta, m} \lambda(\theta, m) u_S(\theta, \sigma^*(m)).$$

Second, suppose  $(\lambda^*, \sigma^*)$  violates R-IC. Then there exists  $\sigma'$  such that

$$\sum_{\theta, m} \lambda^*(\theta, m) u_R(\theta, \sigma'(m)) > \sum_{\theta, m} \lambda^*(\theta, m) u_R(\theta, \sigma^*(m))$$

By continuity, there exist  $\eta > 0$  and  $\varepsilon_4 > 0$  such that for all  $\lambda'$  satisfying  $|\lambda^* - \lambda'| < \varepsilon_4$ , we have

$$\sum_{\theta, m} \lambda'(\theta, m) u_R(\theta, \sigma'(m)) - \sum_{\theta, m} \lambda'(\theta, m) u_R(\theta, \sigma^*(m)) \geq \eta > 0$$

Let  $d \equiv \max_{\theta, a} u_R(\theta, a) - \min_{\theta, a} u_R(\theta, a)$  denote the gap between the Receiver's highest

and lowest payoffs. Let  $\varepsilon_5 \leq \frac{\eta}{d+\eta}$  and  $\varepsilon = \min\{\varepsilon_3, \varepsilon_4, \varepsilon_5\}$ . By assumption, there exists a credible and R-IC finite-sample profile  $(\nu^N, \phi^N, \sigma^N)$  such that  $Pr(|\phi^N(F_\Theta^N) - \lambda^*| \leq \varepsilon) > 1 - \varepsilon$ , and  $\sigma^N = \sigma^*$ . We will show that in the finite sample profile  $(\nu^N, \phi^N, \sigma^N)$ , the Receiver can profitably deviate from  $\sigma^N = \sigma^*$  to  $\sigma'$ , which contradicts  $(\nu^N, \phi^N, \sigma^N)$  being R-IC.

By choosing  $\sigma^*$  the Receiver obtains payoff

$$\sum_{F_\Theta^N \in \mathcal{F}_\Theta^N} P^N(F_\Theta^N) \sum_{\theta, m} \phi^N(\theta, m | F_\Theta^N) u_R(\theta, \sigma^*(m))$$

By contrast, the Receiver obtains

$$\sum_{F_\Theta^N \in \mathcal{F}_\Theta^N} P^N(F_\Theta^N) \sum_{\theta, m} \phi^N(\theta, m | F_\Theta^N) u_R(\theta, \sigma'(m))$$

from choosing  $\sigma'$ . Denote  $E^N \equiv \{F_\Theta^N : |\phi^N(F_\Theta^N) - \lambda^*| \leq \delta\}$  so  $Pr(E^N) > 1 - \varepsilon$ . By switching from  $\sigma^*$  to  $\sigma'$ , the Receiver obtains an extra payoff of

$$\begin{aligned} & \sum_{F_\Theta^N \in \mathcal{F}_\Theta^N} P^N(F_\Theta^N) \left[ \sum_{\theta, m} \phi^N(\theta, m | F_\Theta^N) u_R(\theta, \sigma^*(m)) - \sum_{\theta, m} \phi^N(\theta, m | F_\Theta^N) u_R(\theta, \sigma'(m)) \right] \\ &= \sum_{F_\Theta^N \in E^N} P^N(F_\Theta^N) \left[ \sum_{\theta, m} \phi^N(\theta, m | F_\Theta^N) u_R(\theta, \sigma^*(m)) - \sum_{\theta, m} \phi^N(\theta, m | F_\Theta^N) u_R(\theta, \sigma'(m)) \right] \\ &+ \sum_{F_\Theta^N \notin E^N} P^N(F_\Theta^N) \left[ \sum_{\theta, m} \phi^N(\theta, m | F_\Theta^N) u_R(\theta, \sigma^*(m)) - \sum_{\theta, m} \phi^N(\theta, m | F_\Theta^N) u_R(\theta, \sigma'(m)) \right] \end{aligned}$$

Note that  $\sum_{\theta, m} \phi^N(\theta, m | F_\Theta^N) u_R(\theta, \sigma^*(m)) - \sum_{\theta, m} \phi^N(\theta, m | F_\Theta^N) u_R(\theta, \sigma'(m)) \geq \eta$  for all  $F_\Theta^N \in E^N$ , while for all  $F_\Theta^N \notin E^N$ ,  $\sum_{\theta, m} \phi^N(\theta, m | F_\Theta^N) u_R(\theta, \sigma^*(m)) - \sum_{\theta, m} \phi^N(\theta, m | F_\Theta^N) u_R(\theta, \sigma'(m)) \geq -d$ . Together they imply,

$$\begin{aligned} & \sum_{F_\Theta^N \in \mathcal{F}_\Theta^N} P^N(F_\Theta^N) \left[ \sum_{\theta, m} \phi^N(\theta, m | F_\Theta^N) u_R(\theta, \sigma^*(m)) - \sum_{\theta, m} \phi^N(\theta, m | F_\Theta^N) u_R(\theta, \sigma'(m)) \right] \\ & \geq \eta P^N(E^N) - d(1 - P^N(E^N)). \end{aligned}$$

Since  $P^N(E^N) > 1 - \varepsilon$ , we have

$$\begin{aligned} & \sum_{F_\Theta^N \in \mathcal{F}_\Theta^N} P^N(F_\Theta^N) \left[ \sum_{\theta, m} \phi^N(\theta, m | F_\Theta^N) u_R(\theta, \sigma^*(m)) - \sum_{\theta, m} \phi^N(\theta, m | F_\Theta^N) u_R(\theta, \sigma'(m)) \right] \\ & > (1 - \varepsilon)\eta - \varepsilon d = \eta - \varepsilon(\eta + d) \geq 0 \end{aligned}$$

This contradicts the R-IC of  $(\nu^N, \phi^N, \sigma^N)$ .

□

*Proof of Proposition 6 statement 2.* For each  $N \geq 1$ , define  $\sigma^N = \sigma^*$ ,  $\nu^N \in \arg \min_{\nu \in \mathcal{F}_M^N} |\lambda_M^* - \nu|$ , and  $\phi^N : \mathcal{F}_\Theta^N \rightarrow \cup_{F_\Theta^N \in \mathcal{F}_\Theta^N} X(F_\Theta^N, \nu^N)$  by

$$\phi(F_\Theta^N) \in \arg \max_{\lambda \in X^N(F_\Theta^N, \nu^N)} \sum_{\theta, m} \lambda(\theta, m) u_S(\theta, \sigma^*(m)).$$

By construction, for every  $N$ ,  $(\nu^N, \phi^N, \sigma^N)$  is credible and  $|\sigma^N - \sigma^*| = 0$ . It remains to show that for every  $\varepsilon > 0$ , there exists large enough  $N$ , such that

1.  $|\nu^N - \lambda_M^*| < \varepsilon$ ;
2.  $P(|\phi^N(F_\Theta^N) - \lambda^*| < \varepsilon) > 1 - \varepsilon$ ;
3.  $(\nu^N, \phi^N, \sigma^N)$  is R-IC.

From the denseness of rational numbers, we know that  $\nu^N \rightarrow \lambda_M^*$  as  $N \rightarrow \infty$  so the first statement follows.

To prove the second statement, note that since  $(\lambda^*, \sigma^*)$  is strictly credible,  $\lambda^*$  is the unique maximizer to

$$\max_{\lambda \in \Lambda(\mu, \nu)} \sum_{\theta, m} \lambda(\theta, m) u_S(\theta, \sigma^*(m)).$$

From Lemma 10, the best response correspondence  $B(\mu, \nu) \equiv \arg \max_{\lambda \in \Lambda(\mu, \nu)} \sum_{\theta, m} \lambda(\theta, m) u_S(\theta, \sigma^*(m_j))$  is upper hemi-continuous. Since  $B(\mu, \nu) = \{\lambda^*\}$  is a singleton, from Lemma 12,  $B$  is continuous at  $(\mu, \nu)$ . Therefore, there exists  $\delta > 0$ , so that for any  $(\mu', \nu')$  such  $|\mu - \mu'| < \delta$  and  $|\nu - \nu'| < \delta$ , we have  $|\lambda' - \lambda^*| < \varepsilon$  for every  $\lambda' \in B(\mu', \nu')$ .

From the Glivenko–Cantelli theorem, for large  $N$ ,  $P(|F_\Theta^N - \mu_0| < \delta) > 1 - \varepsilon$ . Pick  $N$  large enough so that  $P(|F_\Theta^N - \mu_0| < \delta) > 1 - \varepsilon$  and  $|\nu^N - \mu_M^*| < \delta$ . Follows from the definition of  $\phi$  and Lemma 11,  $\phi(F_\Theta^N) \in \arg \max_{\lambda \in X^N(F_\Theta^N, \nu^N)} \sum_{\theta, m} \lambda(\theta, m) u_S(\theta, \sigma^*(m)) \subset \arg \max_{\lambda \in \Lambda(F_\Theta^N, \nu^N)} \sum_{\theta, m} \lambda(\theta, m) u_S(\theta, \sigma^*(m))$ . So with at least  $1 - \varepsilon$  probability,  $|\phi(F_\Theta^N) - \lambda^*| < \varepsilon$ .

Lastly, we show that  $(\nu^N, \phi^N, \sigma^N)$  is R-IC for large  $N$ . Since  $(\lambda^*, \sigma^*)$  is strictly R-IC, for any  $\sigma \neq \sigma^*$ ,

$$\sum_{\theta, m} \lambda^*(\theta, m) u_R(\theta, \sigma^*(m)) > \sum_{\theta, m} \lambda^*(\theta, m) u_R(\theta, \sigma(m)).$$

From continuity, there exists  $\eta > 0$  such that for any  $\lambda$  such that  $|\lambda^* - \lambda| < \varepsilon$ ,

$$\sum_{\theta, m} \lambda(\theta, m) u_R(\theta, \sigma^*(m)) - \sum_{\theta, m} \lambda(\theta, m) u_R(\theta, \sigma(m)) \geq \eta > 0.$$

As we have shown, for any  $\varepsilon > 0$ , for large enough  $N$ ,  $Pr(|\phi(F_\Theta^N) - \lambda^*| \leq \varepsilon) \geq 1 - \varepsilon$ . Pick  $\varepsilon \leq \frac{\eta}{d+\eta}$ , then follow from the same argument above, we have

$$\sum_{F_\Theta^N \in \mathcal{F}_\Theta^N} P^N(F_\Theta^N) \sum_{\theta, m} \phi(\theta, m | F_\Theta^N) u_R(\theta, \sigma^*(m)) > \sum_{F_\Theta^N \in \mathcal{F}_\Theta^N} P^N(F_\Theta^N) \sum_{\theta, m} \phi(\theta, m | F_\Theta^N) u_R(\theta, \sigma(m))$$

for any  $\sigma \neq \sigma^*$ . So  $(\nu^N, \phi^N, \sigma^N)$  is R-IC.

□

## B Supplementary Appendix

### B.1 Extension to Infinite Spaces

Suppose  $\Theta$  and  $A$  are compact Polish spaces, and let  $M$  be a Polish space containing  $A$ . An information structure  $\lambda \in \Delta(\Theta \times M)$  is a Borel probability measure on  $\Theta \times M$ . A strategy  $\sigma : M \rightarrow A$  is a measurable function from  $M$  to  $A$ .

An outcome distribution is a Borel measure  $\pi \in \Delta(\Theta \times A)$ . The outcome distribution  $\pi$  is induced by the profile  $(\lambda, \sigma)$  if  $\pi$  is the pushforward measure of  $\lambda$  obtained from the function  $\rho : (\theta, m) \rightarrow (\theta, \sigma(m))$ . That is, for any  $S \in \mathcal{B}(\Theta \times A)$ ,  $\pi(S) = \lambda(\rho^{-1}(S))$ .

**Definition 1\***. A profile  $(\lambda, \sigma)$  is credible if

$$\lambda \in \arg \max_{\lambda' \in D(\lambda)} \int u_S(\theta, \sigma(m)) d\lambda',$$

where  $D(\lambda) = \{\lambda' \in \Delta(\Theta \times M) \mid \lambda'_\Theta = \mu_0, \lambda'_M = \lambda_M\}$ .

**Definition 2\***. A profile  $(\lambda, \sigma)$  is R-IC if for any Receiver's strategy  $\sigma'$ , we have

$$\int u_R(\theta, \sigma(m)) d\lambda \geq \int u_R(\theta, \sigma'(m)) d\lambda.$$

**Definition 3\***. An outcome distribution is stable if it can be induced by a profile that is both credible and R-IC.

#### B.1.1 Extension of Theorem 1

Let  $\{\pi(\cdot|a)\}_{a \in A} \subseteq \Delta(\Theta)$  denote a system of regular conditional probabilities obtained from disintegrating  $\pi$  with  $\pi_A$  (see, for example, [Chang and Pollard, 1997](#)). The following result is an extension of [Theorem 1](#).

**Theorem 1\*.** An outcome distribution  $\pi \in \Delta(\Theta \times A)$  is stable if and only if there exists a Borel set  $E^\circ \subseteq \Theta \times A$  with  $\pi(E^\circ) = 1$  such that

1.  $\pi$  is  $u_S$ -cyclically monotone on  $E^\circ$ : for any sequence  $(\theta^1, a^1), \dots, (\theta^n, a^n) \in E^\circ$  and  $a^{n+1} \equiv a^1$ ,

$$\sum_{i=1}^n u_S(\theta^i, a^i) \geq \sum_{i=1}^n u_S(\theta^i, a^{i+1}).$$

2.  $\pi$  is  $u_R$ -obedient on  $E^\circ$ : for each  $a \in A^\circ \equiv \text{proj}_A(E^\circ)$ , let  $E_a^\circ \equiv \{\theta : (\theta, a) \in E^\circ\}$ , then

$$\int_{E_a^\circ} u_R(\theta, a) d\pi(\theta|a) \geq \int_{E_a^\circ} u_R(\theta, a') d\pi(\theta|a)$$

for all  $a \in A^\circ$  and all  $a' \in A$ .

*Proof.* The “only if” direction: Suppose that  $\pi$  is induced by a credible and R-IC profile  $(\lambda, \sigma)$ . The profile  $(\lambda, \sigma)$  being credible implies

$$\lambda \in \arg \max_{\lambda' \in D(\lambda)} \int u_S(\theta, \sigma(\theta, m)) d\lambda'(\theta, m)$$

Let  $\tilde{u}(\theta, m) \equiv u_S(\theta, \sigma(m))$ . Since  $\tilde{u}(\theta, m)$  is Borel measurable and  $|\tilde{u}(\theta, m)| < \infty$ , by [Beiglböck et al. \(2009\)](#),  $\lambda$  is  $\tilde{u}$ -cyclically monotone: i.e. there exists a Borel set  $F \subseteq \Theta \times M$  such that  $\lambda(F) = 1$  and for every sequence  $(\theta_1, m_1), \dots, (\theta_n, m_n) \in F$ ,

$$\sum_{i=1}^n u_S(\theta_i, \sigma(m_i)) \geq \sum_{i=1}^n u_S(\theta_i, \sigma(m_{i+1})).$$

Consider the function  $\rho : (\theta, m) \rightarrow (\theta, \sigma(m))$ , and define  $E \equiv \rho(F)$ . Since  $\lambda(F) = 1$  and  $\pi$  is the pushforward measure of  $\lambda$  obtained from  $\rho$ , it follows that  $\pi(E) = 1$ . In addition, for any sequence  $(\theta_1, a_1), \dots, (\theta_n, a_n) \in E$ , there exists sequence  $(\theta_1, m_1), \dots, (\theta_n, m_n) \in F$  such that  $a_i = \sigma(m_i)$ . So

$$\sum_{i=1}^n u_S(\theta_i, a_i) = \sum_{i=1}^n u_S(\theta_i, \sigma(m_i)) \geq \sum_{i=1}^n u_S(\theta_i, \sigma(m_{i+1})) = \sum_{i=1}^n u_S(\theta_i, a_{i+1}),$$

which implies that  $\pi$  is  $u_S$ -cyclically monotone on the set  $E$ .

Now for each  $a \in A$ , let  $E_a \equiv \{\theta : (\theta, a) \in E\}$ . Note that  $\pi(E_a|a) = 1$  for  $\pi_A$ -almost all  $a \in A$ , since otherwise there exists  $\tilde{A} \subseteq A$  with  $\pi_A(\tilde{A}) > 0$ , such that for all  $a \in \tilde{A}$ ,

$\pi(E_a|a) < 1$ . This would then imply

$$\begin{aligned}
\pi(\Theta \times \tilde{A}) &= \pi(E \cap (\Theta \times \tilde{A})) \\
&= \int_A \left[ \int_{\Theta} \mathbb{1}_E \times \mathbb{1}_{\Theta \times \tilde{A}} d\pi(\theta|a) \right] d\pi_A(a) \\
&< \int_{\tilde{A}} 1 d\pi_A(a) \\
&= \pi(\Theta \times \tilde{A}),
\end{aligned}$$

which is a contradiction. So  $\pi(E_a|a) = 1$  for  $\pi_A$ -almost all  $a \in A$ . As a result, for all measurable functions  $\phi : \Theta \rightarrow \mathbb{R}$  and all  $a \in A$ , we have

$$\int_{\Theta} g d\pi(\theta|a) = \int_{E_a} g d\pi(\theta|a).$$

Next we establish that for  $\pi_A$ -almost all  $a \in A$ ,

$$\int_{E_a} u_R(\theta, a) d\pi(\theta|a) \geq \int_{E_a} u_R(\theta, a') d\pi(\theta|a) \quad (27)$$

for all  $a' \in A$ . We prove this by proving its contraposition: suppose this is not true, we will show that this implies  $(\lambda, \sigma)$  is not R-IC. Specifically, if (27) does not hold for  $\pi_A$ -almost all  $a \in A$  and all  $a' \in A$ , then there exists  $\hat{A} \in \mathcal{B}(A)$  with  $\pi_A(\hat{A}) > 0$ , and for each  $a \in \hat{A}$ , we can find  $d(a) \in A$  that satisfies

$$\int_{E_a} u_R(\theta, d(a)) d\pi(\theta|a) > \int_{E_a} u_R(\theta, a) d\pi(\theta|a).$$

Since  $u_R(\theta, a)$  is a bounded Carathéodory function, the function

$$g(a, a') \equiv \int_{E_a} u_R(\theta, a') d\pi(\theta|a)$$

is measurable in  $a$  and continuous in  $a'$ , and therefore also Carathéodory. For each  $a \in \hat{A}$ , let  $\phi(a) \equiv \arg \max_{a' \in A} g(a, a')$  denote the maximizers of the Receiver's interim expected payoff. Since  $A$  is compact, by the Measurable Maximum Theorem (see, for example, Theorem 18.19 in [Aliprantis and Border, 2006](#)), the correspondence  $\phi(a)$  admits a measurable selection  $d^* : \hat{A} \rightarrow A$ , such that for all  $a \in \hat{A}$ ,

$$\int_{E_a} u_R(\theta, d^*(a)) d\pi(\theta|a) \geq \int_{E_a} u_R(\theta, d(a)) d\pi(\theta|a) > \int_{E_a} u_R(\theta, a) d\pi(\theta|a).$$

Now define  $f^* = f$  for  $a \in \hat{A}$  and  $f^* = I$  for  $a \notin \hat{A}$ . Clearly  $f^* : A \rightarrow A$  is measurable. In addition,

$$\int_{E_a} u_R(\theta, f^*(a)) d\pi(\theta|a) > \int_{E_a} u_R(\theta, a) d\pi(\theta|a).$$

for all  $a \in \hat{A}$ . Since  $\pi_A(\hat{A}) > 0$ , we have that

$$\begin{aligned} \int_{\Theta \times A} u_R(\theta, f^*(a)) d\pi(\theta, a) &= \int_A \left[ \int_{\Theta} u_R(\theta, f^*(a)) d\pi(\theta|a) \right] d\pi_A(a) \\ &= \int_A \left[ \int_{E_a} u_R(\theta, f^*(a)) d\pi(\theta|a) \right] d\pi_A(a) \\ &> \int_A \left[ \int_{E_a} u_R(\theta, a) d\pi(\theta|a) \right] d\pi_A(a) \\ &= \int_A \left[ \int_{\Theta} u_R(\theta, a) d\pi(\theta|a) \right] d\pi_A(a) \\ &= \int_{\Theta \times A} u_R(\theta, a) d\pi(\theta, a), \end{aligned} \tag{28}$$

Now since  $\pi$  is the pushforward measure of  $\lambda$ , we have

$$\int_{\Theta \times A} u_R(\theta, a) d\pi(\theta, a) = \int_{\Theta \times M} u_R(\theta, \sigma(m)) d\lambda(\theta, m). \tag{29}$$

In addition, let  $\sigma' \equiv f^* \circ \sigma$ , then  $\sigma' : M \rightarrow \mathbb{R}$  is a Borel measurable function on  $M$ , and

$$\begin{aligned} \int_{\Theta \times A} u_R(\theta, f^*(a)) d\pi(\theta, a) &= \int_{\Theta \times M} u_R(\theta, f^* \circ \sigma(m)) d\lambda(\theta, m) \\ &= \int_{\Theta \times M} u_R(\theta, \sigma'(m)) d\lambda(\theta, m). \end{aligned} \tag{30}$$

Plugging (29) and (30) into (28), we have

$$\int_{\Theta \times M} u_R(\theta, \sigma'(m)) d\lambda(\theta, m) > \int_{\Theta \times A} u_R(\theta, \sigma(m)) d\lambda(\theta, m),$$

which is a contradiction to  $(\lambda, \sigma)$  being R-IC. So there exists  $\bar{A} \subseteq A$  with  $\pi_A(\bar{A}) = 1$ , such that

$$\int_{E_a} u_R(\theta, a) d\pi(\theta|a) \geq \int_{E_a} u_R(\theta, a') d\pi(\theta|a)$$

for all  $a \in \bar{A}$  and all  $a' \in A$ .

Define  $E^\circ \equiv E \cap (\Theta \times \bar{A})$ . Note that  $\pi(E^\circ) = 1$ , and  $\pi$  is  $u_R$ -obedient on  $E^\circ$ . In addition, since  $\pi$  is  $u_S$ -cyclically monotone on  $E$  and  $E^\circ \subset E$ , we have that  $\pi$  is  $u_S$ -cyclically monotone on  $E^\circ$ . This completes the proof of the “only if” direction.



*The “if” direction:* Suppose there exists a Borel set  $E^\circ \subseteq \Theta \times A$  with  $\pi(E^\circ) = 1$ , where the outcome distribution  $\pi \in \Delta(\Theta \times A)$  is both  $u_S$ -cyclical monotone and  $u_R$ -obedient. Let the message space  $M = A$ , and consider the profile  $(\lambda, \sigma)$  where  $\lambda \equiv \pi$  and  $\sigma$  is the identity mapping. Clearly,  $(\lambda, \sigma)$  induces  $\pi$ . We will show that  $(\lambda, \sigma)$  is both credible and R-IC.

To see that  $(\lambda, \sigma)$  is R-IC, first note that following a similar argument as the one in the “only if” direction, we have

$$\int_{\Theta} g d\pi(\theta|a) = \int_{E_a^\circ} g d\pi(\theta|a).$$

for all measurable functions  $\phi : \Theta \rightarrow \mathbb{R}$  and  $\pi_A$ -almost all  $a \in A$ . So for all  $\sigma' : A \rightarrow A$ ,

$$\begin{aligned} \int_{\Theta \times A} u_R(\theta, a) d\pi(\theta, a) &= \int_A \int_{\Theta} u_R(\theta, a) d\pi(\theta|a) d\pi_A(a) \\ &= \int_A \int_{E_a^\circ} u_R(\theta, a) d\pi(\theta|a) d\pi_A(a) \\ &\geq \int_A \int_{E_a^\circ} u_R(\theta, \sigma'(a)) d\pi(\theta|a) d\pi_A(a) \\ &= \int_{\Theta \times A} u_R(\theta, \sigma'(a)) d\pi(\theta, a), \end{aligned}$$

which implies that  $(\lambda, \sigma)$  is R-IC.

Next we show  $(\lambda, \sigma)$  is credible. Since  $\pi$  is  $u_S$ -cyclically monotone on  $E^\circ$ , every sequence  $(\theta_1, a_1), \dots, (\theta_n, a_n) \in E^\circ$  satisfies

$$\sum_{i=1}^n u_S(\theta_i, a_i) \geq \sum_{i=1}^n u_S(\theta_i, a_{i+1}),$$

where  $a_{n+1} \equiv a_1$ . Since  $\lambda = \pi$  and  $\sigma$  is the identity mapping, this is equivalent to the existence of a set  $F \subseteq \Theta \times M$  with  $\lambda(F) = 1$ , such that

$$\sum_{i=1}^n u_S(\theta_i, \sigma(m_i)) \geq \sum_{i=1}^n u_S(\theta_i, \sigma(m_{i+1}));$$

for every sequence  $(\theta_1, m_1), \dots, (\theta_n, m_n) \in F$  with  $m_{n+1} = m_1$ . By [Beiglböck et al. \(2009\)](#),  $\lambda$  satisfies

$$\lambda \in \arg \max_{\lambda' \in D(\lambda)} \int_{\Theta \times M} u_S(\theta, \sigma(m)) d\lambda'$$

which means  $(\lambda, \sigma)$  is credible. □

### B.1.2 Extension of Proposition 1

Next, we extend Proposition 1 to infinite spaces. Let  $A$ ,  $\Theta$ , and  $M$  be compact subsets of  $\mathbb{R}$ , and  $A \subseteq M$ .

**Definition 4.** An outcome distribution  $\pi \in \Delta(\Theta \times A)$  is a no-information outcome if there exists  $\hat{a} \in A$  such that  $\pi(\Theta \times \{\hat{a}\}) = 1$ .

For each pair of actions  $a, a' \in A$ , let  $\Theta_0(a, a') \equiv \{\theta : u_R(\theta, a) = u_R(\theta, a')\}$  denote the set of states under which the Receiver is indifferent between  $a$  and  $a'$ .

**Proposition 1\*.** Suppose  $\mu_0(\Theta_0(a, a')) < 1$  for all distinct  $a, a' \in A$ . In addition, suppose  $u_S : \Theta \times A \rightarrow \mathbb{R}$  is strictly supermodular and  $u_R : \Theta \times A \rightarrow \mathbb{R}$  is submodular, then any stable outcome distribution must be a no-information outcome.

*Proof.* Let  $\pi$  be a stable outcome distribution. By Theorem 1 and Lemma 1, there exists a Borel set  $E^\circ \subseteq \Theta \times A$  with  $\pi(E^\circ) = 1$ , such that

1.  $\pi$  is comonotone on  $E^\circ$ : for all  $(\theta, a), (\theta', a') \in E$ ,  $a < a'$  implies  $\theta \leq \theta'$ ; and
2.  $\pi$  is  $u_R$ -obedient on  $E^\circ$ : for each  $a \in A^\circ \equiv \text{proj}_A(E^\circ)$ , let  $I_a \equiv \{\theta : (\theta, a) \in E^\circ\}$ , then

$$\int_{I_a} u_R(\theta, a) d\pi(\theta|a) \geq \int_{I_a} u_R(\theta, a') d\pi(\theta|a)$$

for all  $a \in A^\circ$  and  $a' \in A$ .

From  $u_R$ -obedience, we know that for all  $a, a' \in A^\circ$  with  $a < a'$ , we have

$$\int_{I_a} [u_R(\theta, a) - u_R(\theta, a')] d\pi(\theta|a) \geq 0 \geq \int_{I_{a'}} [u_R(\theta, a) - u_R(\theta, a')] d\pi(\theta|a'). \quad (31)$$

In addition, comonotonicity implies that  $\theta \leq \theta'$  for all  $\theta \in I_a, \theta' \in I_{a'}$ . Since  $u_R$  is submodular, we have  $u_R(\theta, a) - u_R(\theta, a') \leq u_R(\theta', a) - u_R(\theta', a')$  for all  $\theta \in I_a$  and  $\theta' \in I_{a'}$ , which implies

$$\sup_{\theta \in I_a} \{u_R(\theta, a) - u_R(\theta, a')\} \leq \inf_{\theta' \in I_{a'}} \{u_R(\theta', a) - u_R(\theta', a')\},$$

and therefore

$$\begin{aligned}
\int_{I_a} [u_R(\theta, a) - u_R(\theta, a')] d\pi(\theta|a) &\leq \int_{I_a} \sup_{\theta \in I_a} [u_R(\theta, a) - u_R(\theta, a')] d\pi(\theta|a) \\
&= \sup_{\theta \in I_a} \{u_R(\theta, a) - u_R(\theta, a')\} \\
&\leq \inf_{\theta' \in I_{a'}} \{u_R(\theta', a) - u_R(\theta', a')\} \\
&\leq \int_{I_{a'}} [u_R(\theta, a) - u_R(\theta, a')] d\pi(\theta|a').
\end{aligned} \tag{32}$$

Combining (31) and (32), we have for all  $a, a' \in A^\circ$ ,  $a < a'$ ,

$$\int_{I_a} [u_R(\theta, a) - u_R(\theta, a')] d\pi(\theta|a) = \sup_{\theta \in I_a} \{u_R(\theta, a) - u_R(\theta, a')\} = 0$$

and

$$\int_{I_{a'}} [u_R(\theta, a) - u_R(\theta, a')] d\pi(\theta|a') = \inf_{\theta' \in I_{a'}} \{u_R(\theta', a) - u_R(\theta', a')\} = 0.$$

This implies that for all  $a, a' \in A^\circ$ ,  $a < a'$ , we have

$$\begin{aligned}
u_R(\theta, a) - u_R(\theta, a') &\leq 0 \text{ for all } \theta \in I_a, \\
\text{with } u_R(\theta, a) &= u_R(\theta, a') \text{ for } \pi(\cdot|a)\text{-almost all } \theta \in I_a;
\end{aligned} \tag{33}$$

and

$$\begin{aligned}
u_R(\theta', a) - u_R(\theta', a') &\geq 0 \text{ for all } \theta' \in I_{a'}, \\
\text{with } u_R(\theta', a) &= u_R(\theta', a') \text{ for } \pi(\cdot|a')\text{-almost all } \theta' \in I_{a'},
\end{aligned} \tag{34}$$

For each  $a \in A^\circ$ , let  $N(a) \equiv \{\theta \in I_a : u_R(\theta, a) \neq u_R(\theta, a') \text{ for some } a' \in A^\circ\}$  denote the set of states in  $I_a$  under which the Receiver is not indifferent towards all actions in  $A^\circ$ . We want to show that  $\pi(N(a)|a) = 0$  for each  $a \in A^\circ$ . Note that this does not follow directly from (33) and (34) since  $A^\circ$  may be uncountably infinite, and an uncountable union of  $\pi(\cdot|a)$ -null sets may no longer be a  $\pi(\cdot|a)$ -null set.

However, note that since  $u_R$  is submodular, for any  $a' > a$ , if  $u_R(\theta, a) - u_R(\theta, a') < 0$ , then  $u_R(\theta', a) - u_R(\theta', a') < 0$  for all  $\theta' < \theta$ . Similarly, for any  $a' < a$ , if  $u_R(\theta, a') - u_R(\theta, a) > 0$ , then  $u_R(\theta', a') - u_R(\theta', a) > 0$  for all  $\theta' > \theta$ . This means that  $N(a)$  is the union of nested sets that are located at either the lower or upper ends of  $I_a$ . We will exploit this structure to show  $\pi(N(a)|a) = 0$ .

For each  $a' \in A^\circ$ ,  $a' > a$ , let us define

$$\hat{N}(a'|a) \equiv \left\{ \theta \in I_a : u_R(\theta, a) - u_R(\theta, a') < 0 \right\},$$

and

$$\hat{\theta}(a'|a) \equiv \sup \left\{ \theta \in I_a : u_R(\theta, a) - u_R(\theta, a') < 0 \right\}.$$

It follows that

$$(-\infty, \hat{\theta}(a'|a)) \cap I_a \subseteq \hat{N}(a'|a) \subseteq (-\infty, \hat{\theta}(a'|a)] \cap I_a \quad (35)$$

and  $\pi(\hat{N}(a'|a)|a) = 0$ .

Analogous, for each  $a' \in A^\circ$ ,  $a' < a$ , define

$$\tilde{N}(a'|a) \equiv \left\{ \theta \in I_a : u_R(\theta, a') - u_R(\theta, a) > 0 \right\},$$

and

$$\tilde{\theta}(a'|a) \equiv \inf \left\{ \theta \in I_a : u_R(\theta, a') - u_R(\theta, a) > 0 \right\},$$

then

$$(\tilde{\theta}(a'|a), \infty) \cap I_a \subseteq \tilde{N}(a'|a) \subseteq [\tilde{\theta}(a'|a), \infty) \cap I_a \quad (36)$$

and  $\pi(\tilde{N}(a'|a)|a) = 0$ .

Let  $\hat{N}(a) \equiv \cup_{a' \in A^\circ, a' > a} \hat{N}(a'|a)$  and  $\tilde{N}(a) \equiv \cup_{a' \in A^\circ, a' < a} \tilde{N}(a'|a)$ , then we have  $N(a) = \hat{N}(a) \cup \tilde{N}(a)$ . In order to show  $\pi(N(a)|a) = 0$ , it suffices to show both  $\pi(\hat{N}(a)|a) = 0$  and  $\pi(\tilde{N}(a)|a) = 0$ . Below we will show  $\pi(\hat{N}(a)|a) = 0$ . The fact that  $\pi(\tilde{N}(a)|a) = 0$  follows from similar arguments.

Let  $\hat{\theta}(a) \equiv \sup_{a' \in A^\circ, a' > a} \hat{\theta}(a'|a)$ . By (35) and the definition of  $\hat{N}(a)$ , we have

$$(-\infty, \hat{\theta}(a)) \cap I_a \subseteq \hat{N}(a) \subseteq (-\infty, \hat{\theta}(a)] \cap I_a.$$

However, note that if  $\hat{\theta}(a) \in \hat{N}(a)$ , then  $\hat{\theta}(a) \in \hat{N}(a'|a)$  for some  $a' \in A^\circ$  with  $a' > a$ , and this would imply  $\pi(\{\hat{\theta}(a)\}|a) = 0$  since  $\pi(\hat{N}(a'|a)|a) = 0$  for all  $a' \in A^\circ$  with  $a' > a$ . Therefore, in order to prove  $\pi(\hat{N}(a)|a) = 0$ , it suffices to prove that  $\pi((-\infty, \hat{\theta}(a)) \cap I_a | a) = 0$ .

To this end, note that  $(-\infty, \hat{\theta}(a)) = \cup_{n=1}^{\infty} (-\infty, \hat{\theta}(a) - 1/n)$ . Since  $\hat{\theta}(a) \equiv \sup_{a' \in A^\circ, a' > a} \hat{\theta}(a'|a)$ , for each  $n \geq 1$ ,  $(-\infty, \hat{\theta}(a) - 1/n) \subseteq (-\infty, \hat{\theta}(a'|a))$  for some  $a' \in A^\circ$  with  $a' > a$ . So for each  $n \geq 1$ ,

$$\pi((-\infty, \hat{\theta}(a) - 1/n) \cap I_a | a) \leq \pi((-\infty, \hat{\theta}(a'|a)) \cap I_a | a) \leq \pi(N(a'|a) | a) = 0.$$

As a result, we have

$$\begin{aligned}\pi\left((-\infty, \hat{\theta}(a)) \cap I_a \mid a\right) &= \pi\left(\bigcup_{n=1}^{\infty} \left((-\infty, \hat{\theta}(a) - 1/n) \cap I_a\right) \mid a\right) \\ &\leq \sum_{n=1}^{\infty} \pi\left((-\infty, \hat{\theta}(a) - 1/n) \cap I_a \mid a\right) = 0,\end{aligned}$$

so  $\pi(\hat{N}(a) \mid a) = 0$ .

Using similar arguments as above, we can establish that  $\pi(\tilde{N}(a) \mid a) = 0$  as well, so  $\pi(N(a) \mid a) = \pi(\hat{N}(a) \cup \tilde{N}(a) \mid a) = 0$ . For each  $a \in A^\circ$ , let

$$\hat{\Theta}_0(a) \equiv \{\theta \in I_a : u_R(\theta, a) = u_R(\theta, a') \text{ for all } a' \in A^\circ\} = [N(a)]^c,$$

so  $\pi(\hat{\Theta}_0(a) \mid a) = 1 - \pi(N(a) \mid a) = 1$ .

Let  $\hat{\Theta}_0 \equiv \{\theta \in \Theta : u_R(\theta, a) = u_R(\theta, a') \text{ for all } a, a' \in A^\circ\}$ . We have

$$\mu_0(\hat{\Theta}_0) = \pi(\hat{\Theta}_0 \times A) = \pi((\hat{\Theta}_0 \times A) \cap E),$$

so

$$\begin{aligned}\mu_0(\hat{\Theta}_0) &= \int_{\Theta \times A} \mathbf{1}_{\hat{\Theta}_0 \times A} \cdot \mathbf{1}_E d\pi(\theta, a) \\ &= \int_A \left[ \int_{\Theta} \mathbf{1}_{\hat{\Theta}_0 \times A} \cdot \mathbf{1}_E d\pi(\theta \mid a) \right] d\pi_A(a) \\ &= \int_A \left[ \int_{I_a} \mathbf{1}_{\hat{\Theta}_0 \times A} d\pi(\theta \mid a) \right] d\pi_A(a) \\ &= \int_A \pi(\hat{\Theta}_0(a) \mid a) d\pi_A(a) \\ &= \int_A 1 d\pi_A(a) = 1.\end{aligned}$$

Recall that  $\Theta_0(a, a') \equiv \{\theta \in \Theta : u_R(\theta, a) = u_R(\theta, a')\}$  and by our assumption  $\mu_0(\Theta_0(a, a')) < 1$  for all distinct  $a, a' \in A$ . Since  $\hat{\Theta}_0 \equiv \{\theta \in \Theta : u_R(\theta, a) = u_R(\theta, a') \text{ for all } a, a' \in A^\circ\}$  and we have established that  $\mu_0(\hat{\Theta}_0) = 1$ ,  $A^\circ$  must be a singleton set. Since  $\pi(\Theta \times A^\circ) \geq \pi(E^\circ) = 1$ , it follows that  $\pi$  is a no-information outcome.  $\square$

## B.2 Credible Persuasion in Games

In this section, we generalize the framework in [Section 2.1](#) to a setting with multiple Receivers, where the Sender can also take actions after information is disclosed. We also allow the state space and action space to be infinite.

Consider an environment with a single Sender (she) and  $r$  Receivers (each of whom is a he). The Sender has action set  $A_S$  while each Receiver  $i \in \{1, \dots, r\}$  has action set  $A_i$ . Let  $A = A_S \times A_1 \times \dots \times A_r$  denote the set of action profiles. Each player has payoff function  $u_i : \Theta \times A \rightarrow \mathbb{R}$ ,  $i = S, 1, \dots, r$ , respectively. The state space  $\Theta$  and action spaces  $A_i$  are Polish spaces endowed with their respective Borel sigma-algebras. Players hold full-support common prior  $\mu_0 \in \Delta(\Theta)$ . We refer to  $G = (\Theta, \mu_0, A_S, u_S, \{A_i\}_{i=1}^r, \{u_i\}_{i=1}^r)$  as the base game.

Let  $M$  be a Polish space that contains  $A$ . The Sender chooses an information structure  $\lambda \in \Delta(\Theta \times M)$  where  $\lambda_\Theta = \mu_0$ : note that this formulation implies that the information structure generates *public* messages observed by all Receivers. Together the information structure and the base game constitute a Bayesian game  $\mathcal{G} = \langle G, \lambda \rangle$ , where:<sup>28</sup>

1. At the beginning of the game a state-message pair  $(\theta, m)$  is drawn from the information structure  $\lambda$ ;
2. The Sender observes  $(\theta, m)$  while the Receivers observe only  $m$ ; and
3. All players choose an action simultaneously.

A strategy profile  $\sigma : \Theta \times M \rightarrow A$  in  $\mathcal{G}$  consists of a Sender's strategy  $\sigma_S : \Theta \times M \rightarrow A_S$  and Receivers' strategies  $\sigma_i : M \rightarrow A_i$ ,  $i = 1, \dots, r$ . For each profile of Sender's information structure and players' strategies  $(\lambda, \sigma)$ , players' expected payoffs are given by

$$U_i(\lambda, \sigma) = \int_{\Theta \times M} u_i(\theta, \sigma(\theta, m)) d\lambda(\theta, m) \quad \text{for } i = S, 1, \dots, r.$$

We now generalize the notion of credibility and incentive compatibility in [Section 2](#) to the current setting. For each  $\lambda$ , let  $D(\lambda) \equiv \{\lambda' \in \Delta(\Theta \times M) : \lambda'_\Theta = \mu_0, \lambda'_M = \lambda_M\}$  denote the set of information structures that induce the same distribution of messages as  $\lambda$ . [Definition 5](#) is analogous to [Definition 1](#), which requires that given the players' strategy profile, no deviation in  $D(\lambda)$  can be profitable for the Sender.

**Definition 5.** A profile  $(\lambda, \sigma)$  is **credible** if

$$\lambda \in \arg \max_{\lambda' \in D(\lambda)} \int u_S(\theta, \sigma(\theta, m)) d\lambda'(\theta, m). \quad (37)$$

In addition, [Definition 6](#) generalizes [Definition 2](#), and requires players' strategies to form a Bayesian Nash equilibrium of the game  $\langle G, \lambda \rangle$ .

---

<sup>28</sup>The information structure  $\lambda$  can be viewed as “additional information” observed by both the Sender and the Receivers, on top of the base information structure where the Sender observes the state and the Receivers do not observe any signal.

**Definition 6.** A profile  $(\lambda, \sigma)$  is *incentive compatible (IC)* if  $\sigma$  is a Bayesian Nash equilibrium in  $\mathcal{G} = \langle G, \lambda \rangle$ . That is,

$$\sigma_S \in \arg \max_{\sigma'_S: \Theta \times M \rightarrow A_S} U_S(\lambda, \sigma'_S, \sigma_{-S}) \quad \text{and} \quad \sigma_i \in \arg \max_{\sigma'_i: M \rightarrow A_i} U_i(\lambda, \sigma'_i, \sigma_{-i}) \quad \text{for } i = 1, \dots, r. \quad (38)$$

Note that in Definition 5, when the Sender deviates to a different information structure, say  $\lambda'$ , we use the original strategy profile  $\sigma(\theta, m)$  to predict players' actions in the ensuing Bayesian game  $\langle G, \lambda' \rangle$ . One might worry that the Sender may simultaneously change not only her information structure but also her strategy  $\sigma_S(\theta, m)$  in  $\langle G, \lambda' \rangle$ . This, however, is unnecessary since the Sender's optimal strategy in  $\langle G, \lambda' \rangle$  will remain unchanged: the Sender knows  $\theta$  perfectly, her best response in  $\langle G, \lambda' \rangle$  depends only on  $\theta$  and the Receivers' actions (and not on her own information structure).

### B.3 An Example

In this section, we provide an example in which the Sender can benefit from credible persuasion, but cannot achieve her optimal full-commitment payoff. This example also corresponds to the first case of Proposition 2.

The prior belief is  $\mu_0$  with  $\mu_0(\theta_H) = 0.4$ . Note that both the Sender's and the Receiver's payoffs are supermodular. The horizontal axis  $\mu$  in the graph represents the probability assigned to  $\theta_H$  by the posterior belief.

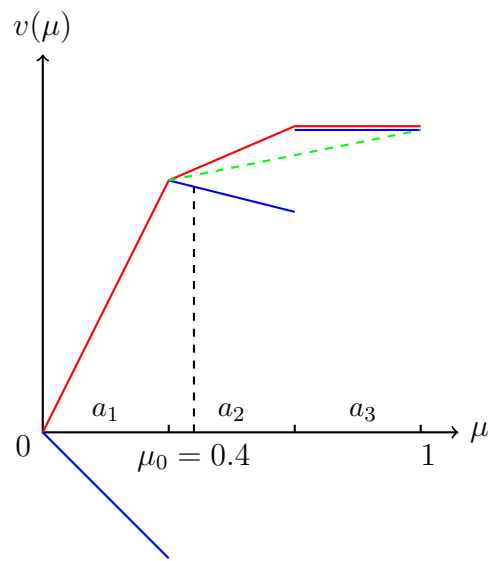
According to the concavification, the optimal full-commitment information structure  $\lambda^*$  induces two posterior beliefs  $\mu = 1/3$  and  $\mu = 2/3$ , with the Receiver's strategy  $\sigma^*$  playing  $a_2$  when  $\mu = 1/3$  and  $a_3$  when  $\mu = 2/3$ . However in this case the support of the outcome distribution is  $\{(\theta_L, a_2), (\theta_L, a_3), (\theta_H, a_2), (\theta_H, a_3)\}$ . This outcome distribution is not comonotone, so  $(\lambda^*, \sigma^*)$  is not credible.

The optimal credible information structure  $\lambda^\circ$  is represented by the green dashed line: it induces two posteriors,  $\mu = 1/3$  and  $\mu = 1$ , with the Receiver strategy  $\sigma^\circ$  playing  $a_2$  when  $\mu = 1/3$  and  $a_3$  when  $\mu = 1$ . In particular, the support of the outcome distribution is  $\{(\theta_L, a_2), (\theta_H, a_2), (\theta_H, a_3)\}$  which is comonotone, so  $(\lambda^\circ, \sigma^\circ)$  is credible.

|              |       |       |       |
|--------------|-------|-------|-------|
| $u_S$        | $a_1$ | $a_2$ | $a_3$ |
| $\theta = H$ | -1    | 0.5   | 0.8   |
| $\theta = L$ | 0     | 0.75  | 0.8   |

|              |       |       |       |
|--------------|-------|-------|-------|
| $u_R$        | $a_1$ | $a_2$ | $a_3$ |
| $\theta = H$ | 0     | 2     | 3     |
| $\theta = L$ | 3     | 2     | 0     |

Sender and Receiver’s payoffs



Concavification