

Perfect Robust Implementation by Private Information Design*

Maxim Ivanov[†]

December 2022

Abstract

This paper studies the general principal-agent framework in which the principal aims to implement his first-best action that is monotone in the unknown state. The principal privately selects a signal structure of the agent whose payoff depends on the principal's action, the state and the privately known type. The agent privately observes the generated signal and reports it to the principal who takes an action. We show that by randomizing between two perfectly informative signal structures, the principal can elicit perfect information from the agent about the state and implement his first-best action regardless of the agent's type. The key idea is that a randomization between signal structures forms posterior beliefs, which induce actions with opposite reactions to agent's messages. As to the economic application, we consider the bilateral-trade model with non-quasilinear preferences of players and show that the seller can extract full surplus from the buyer with private multi-dimensional information.

JEL classification: C72, D81, D82, D83

Keywords: information design; Bayesian persuasion; mechanism design, surplus extraction.

1 Introduction

This paper studies the benefits of private information design as a novel implementation tool in economic environments. The term *private* refers to a situation in which the information

*I am grateful to Kalyan Chatterjee, Mikhail Drugov, Seungjin Han, Peicong Hu, Vijay Krishna, Ming Li, Wooyoung Lim, Elliot Lipnowski, Tatiana Mayskaya, Alan Miller, Alessandro Pavan, Joel Sobel, Chen Zhao, Charles Zheng, and the audiences at Concordia University, Western University, Canadian Economic Theory Conference 2022, Canadian Economic Association Meetings 2022, Stony Brook International Conference on Game Theory 2022, EEA-ESEM 2022 Conference, Econometric Society European Winter Meeting 2022, FES-ICEF-NES seminar, and CUHK-HKU-HKUST Joint Theory Seminar for helpful comments. A part of this work was written while I visited the University of Toronto. I especially thank Heski Bar-Isaac for hospitality and insightful discussions. All errors are mine.

[†]Department of Economics, McMaster University, 1280 Main Street Hamilton, ON, Canada L8S 4M4. Phone: (905) 525-9140. Fax: (905) 521-8232. Email: mivanov@mcmaster.ca.

technology, which provides signals to players about an unknown variable and is a choice of the information designer, is privately known to the designer only. However, the signals generated by the technology are remain private knowledge of the addressees. In this paper, we apply the private information design to the general principal-agent framework. As the main result, we show that the principal (he) can elicit perfect information from the agent (she) and implement his first-best outcome in a simple way by privately designing the agent's information structure. This result holds even if the principal's preferences are incompletely defined, the agent's preferences do or do not depend on this information and/or the agent has separate private information about her preferences.

Before discussing the results, we start with a description of the environment. Its generality stems from substantially relaxing the standard assumptions about players' preferences and their information compared to standard models. First, our model uses only the principal's *ideal action*, which is a monotone function of the unknown *state*. In other words, our results hold for all principal's preferences with a given ideal action function. This type of robustness is economically important, because learning about first-best choices of economic agents—for example, by observing their behavior in the past—is often much easier comparing to obtaining information about non-preferred alternatives. Second, the agent's preferences are represented by any function which is monotone (and differentiable) in the principal's action. Third, the agent's payoff depends on multi-dimensional information. Besides the state, it is also represented by the agent's *type*. The state affects the first-best decision of the principal and potentially the payoff of the agent. For example, in models of trade the state can reflect the product quality, which determines the valuations of the object by both the seller (the principal) and the buyer (the agent). The agent's type reflects only the characteristics of the agent's payoff function, such as the measure of her risk-aversion. As a special case, the agent's preferences may be independent of the state and/or the type. Another key difference, which is also the third factor, is that the agent is a priori privately informed about her type, but uninformed about the state. The uncertainty about the state is a key motive in communication between the agent and the principal.

In short, our framework substantially extends the setup in the related paper by Ivanov and Sam (2022) who study a similar information design problem in the Crawford and Sobel's (1982) cheap-talk framework, in three dimensions. First, they assume that the agent is ex-ante uninformed, while she is privately informed about her type in our setup. Second, the agent in their model must have a unique ideal action in each state. Third, the agent's payoff must be strictly supermodular in the action and the state. (The second and the third assumptions are necessary in most cheap-talk models.) An implication of the last assumption is the dependence of the agent's ideal action and, hence, her payoff, on the state. In contrast, neither of these assumptions about the agent's payoff—the existence of the ideal action, the dependence of this action on the state, or the strict supermodularity—is imposed in our setup. All of them are replaced with the monotonicity of the agent's payoff in the principal's action. As a result, our framework covers a substantially broader class of economic applications than their. For example, our leading application—bilateral trade—does not fit into their model. Specifically, since the buyer's payoff is decreasing in her payment, she does not have an ideal action for any valuation of the product (i.e., the state). Furthermore, some

payoff functions violate the supermodularity condition as well.¹ It is also worth noting that monotonic preferences of the agent represent the worst-case scenario for the principal in terms of information extraction possibilities. The reason is that they precludes any meaningful communication in the model with an arbitrary publicly known signal structure of the agent.

We now turn to the information design part of the model. The information about the state is obtained by the agent whose information technology, called a *signal structure* hereafter, is privately selected by the principal.² Specifically, the principal randomizes between publicly known signal structures, where the realized signal structure is known to the principal only. Once the principal assigns this signal structure to the agent (without informing her about it), the structure generates a signal about the state. Upon privately observing the signal, the agent sends a report to the principal who takes an action based on the report and the realized signal structure.

The first main result of the paper establishes that the principal can elicit perfect information about all states and implement his first-best decision regardless of the agent's type. Moreover, it can be achieved in a simple way by randomizing between two deterministic and perfectly informative signal structures (called *signal functions*). That is, each signal function maps the state into a single signal, and knowing this function allows one to perfectly infer the state from the agent's signal. The key idea behind this result is that a randomization between signal structures creates the uncertainty for the agent about the impact of her message on the principal's action. The principal exploits this uncertainty by selecting signal functions with the 'opposite monotonicities' in the state: a signal generated by one signal function is increasing in the state, while the signal generated by the other signal function is decreasing. This implies opposite reactions of posterior states (induced by each signal function) to the signal, and, hence, opposite reactions of the principal's decisions to the agent's message. That is, any distortion of the signal in an attempt to increase the agent's payoff under one signal function is offset by the marginal losses under the other signal function. By properly selecting signal functions, the principal can sustain agent's truth-telling. Finally, because the principal knows the realized signal function, he can infer the state from the agent's report and implement his first-best action. Furthermore, this result is robust to the agent's privately known type under the local separability condition on the agent's payoff. An implication of this condition is that the agent's marginal benefits from lie and, hence, the optimal signal structures, become invariant to the agent's type.

Potential applications of our results can be illustrated by the following example. Consider a large organization (a firm, a public institution, etc.) with a standard vertical organizational structure: the upper leadership, the middle management, and regular employees. The upper leadership wants to collect the decision-relevant information from employees about some parameter, say, the human resource (HR) effectiveness, by conducting an ad-hoc survey. The employees' responses, however, are collected, processed, and reported by the HR manager herself whose benefits—reputation, bonuses, career opportunities—are increasing in this parameter. The HR manager thus has the incentive to misreport the acquired information.³

¹For example, if the buyer's payoff is a convex function of the difference between her product valuation and the payment, then it is submodular and, thus, violates the strict supermodularity condition.

²The higher opportunity costs of the principal compared to those of the agent is a common reason for assigning the task of collecting information to the agent rather than acquiring it directly.

³In a survey by the consulting firm McKinsey & Company (2007), 36% of top executives responded that

In order to preclude such manipulations, the upper leadership can randomize between two surveys whose questions are known to the leadership, but not to the HR manager.⁴ The key idea is that the survey questions are formulated quantitatively, where the relationship between the reported number and the parameter—that is, an increasing or a decreasing signal function—is known only to the upper leadership. One survey, for example, might ask respondents to rate the HR’s *effectiveness* on a scale of 0 to 100, which corresponds to an increasing signal function. The other survey would ask to evaluate the *ineffectiveness*, with a score of 100 indicating a completely unproductive or mismanaged HR team. Because the HR manager remains uncertain about the actual survey upon receiving employees’ responses, then misreporting would distort the signal in the ‘wrong direction’ with a positive probability and thus stochastically penalize her.⁵

As the leading application, we consider the generalized bilateral-trade model with non-quasilinear preferences of players and private multidimensional information of the buyer. As typically assumed in the literature (Bergemann and Pesendorfer, 2007; Li and Shi, 2019; and Bergemann et al., 2022), the seller determines: i) the terms of trade, i.e., a mechanism, which enforces an allocation and payments on the basis of the buyer’s report; and ii) the buyer’s signal structure.⁶ In our model, the buyer’s information is represented by two variables: the state and the type. The state reflects the intrinsic characteristics of the object (the quality), which determines the buyer’s willingness to pay and can also affect the seller’s payoff. The buyer’s type reflects her preferences only, for instance, the degree of the risk-aversion. Also, the buyer is a priori uninformed about the state and informed about her type. Upon observing a private signal generated by the signal structure and learning the type, the buyer sends a message to the mechanism, which enforces the terms of trade. Similarly to the main framework, the outcome of the mechanism depends on both the buyer’s message and the privately known signal structure. This creates the uncertainty for the buyer regarding the impact of her message on the mechanism.

Our second main result shows that the seller can extract the full surplus from the buyer upon eliciting the perfect information about the state. The full information and surplus extraction is feasible by employing private signal structures similar to those in the general model. The only difference is that a randomization between signal functions creates the buyer’s uncertainty about the *target* subinterval of states rather than the entire state space. The target states reflect the fact that the seller might be interested to sell the object only if the buyer’s willingness to pay—the highest payment, which make her indifferent between trading and taking an outside option—is high enough so that it exceeds

managers hide, restrict, or misrepresent information at least “somewhat” frequently.

⁴It is also assumed that the manager cannot access the questions in another way, for example, by spying on employees or cooperating with them.

⁵Two key factors differentiate this example from university processes, when students’ feedback at the end of each semester is going directly to academic services by bypassing professors. First, the ad hoc character of the survey leaves the manager uncertain about the questions in it, while standard questions in the students’ evaluation forms do not change over time and are thus publicly known. Second, involving external evaluators (similar to academic services at universities) for an ad-hoc survey would require substantial costs per survey, which exceed the costs of running routine surveys at universities.

⁶In practice, sellers often allow buyers to try or test the product before purchasing, provide a demo version of the product, or let buyers gather additional information about products in order to assess their quality.

the seller’s benefits from keeping the object.⁷ Specifically, the optimal private signal structure randomizes between two signal functions with the opposite reactions of signals to states. Also, the seller sells the product if and only if the posterior state induced by the message belongs to the target subset and charges the buyer with her willingness to pay in each posterior state. Then, the opposite monotonicities of signal functions in the state imply the opposite monotonicities of the buyer’s payments in her message under these signal functions. That is, any marginal benefits from distorting the observed signal in an attempt to reduce the payment under one signal function are offset by the higher payment under the other signal function. This trade-off sustains the incentive-compatibility of the mechanism. Importantly, it does not depend on the absolute value of the buyer’s willingness to pay. As a result, the seller can charge the buyer with her willingness to pay (which does not depend on the buyer’s type) upon learning the state.

In this light, our paper extends the related models in the mechanism design literature in three dimensions. First, it does not assume that the players’ preferences are quasilinear. It is a conceptual extension. If the parties’ preferences are quasilinear in the state, and the seller’s payoff from the product is constant, then he can extract the full surplus by informing the buyer whether her valuation of the object is above or below that of the seller and setting the price at the higher posterior value (Saak, 2006). This construction, however, cannot be extended to setups in which the players’ preferences take a more general form. For example, if the buyer is risk-averse, then hiding information about the product valuation reduces her willingness to pay and, as a consequence, does not allow the seller to extract the full surplus. In addition, if the seller’s payoff also depends on the buyer’s type, then he must also learn this information before making a decision about selling the object. Otherwise, there is a chance of selling the object whose value to the seller is above the price. This results in the allocative inefficiency, which does not allow the seller to extract all gains from trade.

Second, we show that the full surplus extraction is robust to the buyer’s privately knowledge of her preferences. As noted by Krämer (2020), it is a hard and generally unsolvable problem if the seller employs the private information design in the model with quasilinear preferences, such that the players’ valuations depend on the ex-ante unknown state and the privately known buyer’s type. Furthermore, it is actually insolvable if the state and type are independent as we assumed in our model. As we show, however, the problem can be circumvented under two conditions on the buyer’s payoff function. First, the buyer’s willingness to pay is solely determined by the state (i.e., the product quality). As we elaborate in detail below, this condition is not novel and imposed, for instance, in models of auctions with non-linear buyers’ preferences.⁸ In this case, the type parameterizes the buyer’s payoff function and reflects, for instance, its non-linearity (e.g., the measure of risk-aversion), which does not affect her willingness to pay for the object. Second, the buyer’s marginal payoff with respect to the payment at the payment equal to the buyer’s willingness to pay can be factorized into separate functions of the state and the type. Notably, both conditions are local as they must hold only for the payment equal to the buyer’s willingness to pay. In the leading example, we show that these conditions hold for *all* payoff functions

⁷For low states, the precision of the buyer’s information is not substantial. In fact, the seller can select identical perfectly informative signal functions for these states. Thus, even though the buyer perfectly infers the state from the signal, she still prefers not to purchase the object.

⁸See, for instance, Section 4.1 in Krishna (2009).

that depend on: i) the difference between the state and the payment; and ii) the type, which determines the shape of the payoff as a function of this difference. As we show, under these conditions the seller is still able to extract the full surplus by using the private signal functions, which are robust to the buyer’s type.

Third, in contrast to most of the literature on full surplus extraction, the set of states (and potentially types) in our model is continuous. It is also not a purely technical extension. Eliciting information, which can take a finitely many values, is generally a less difficult problem, since the set of the buyers’ incentive-compatibility constraints is substantially smaller. Importantly, buyers cannot distort their information locally by mimicking nearby values. At the same time, local incentive-compatibility plays a crucial role in the mechanism design (Myerson, 1981). In addition, there are qualitative differences between spaces of finite and continuous distributions. Because of them, the results by Crémer and McLean (1988) and Krähmer (2020) about full surplus extraction from buyers with discrete valuations—which rely on eliciting posterior distributions from other correlated variables—are generically impossible for continuous distributions of states (Heifetz and Neeman, 2006).⁹

Literature

Besides the mentioned work by Ivanov and Sam (2022), who first introduced the idea of randomization between signal functions with the opposite monotonicities to extract information from the agent in the cheap-talk model, another related work is by Watson (1996). He showed that the principal can elicit information from the perfectly informed agent and implement her first-best outcome by maximally confusing the agent about the relationship between his message and the principal’s decision, where the relationship is the private information of the principal. It is achieved by (approximately) uniform randomization over the set of *all* bijective mappings from reported states to principal’s ideal actions. In contrast, our construction is based on keeping the agent uncertain about the state, which *leads to* uncertainty about the relationship between his message and the principal’s action. Next, our construction is very simple and involves randomizing between only two signal functions, i.e., bijective mappings. In other words, its main goal is to keep the agent minimally uncertain about the relationship between the message and the action rather than maximally uncertain as in Watson’s model. Finally, we show that the shapes of signal functions play a critical role in sustaining agent’s truth-telling.¹⁰

In the context of the mechanism design, the most relevant paper to ours is Krähmer (2020) who first used the private information design to demonstrate the possibility of the full surplus extraction in the bilateral-trade model. There are several key distinctions between our and Krähmer’s papers. First, he considers the quasilinear preferences of the buyer, whereas the buyers’ preferences in our model are of a general form. Second, while the Krähmer’s general model considers a privately informed buyer, the full-surplus extraction is demonstrated for an a priori uninformed buyer only. We establish this result even if the

⁹However, almost full surplus extraction can be attained by partitioning the state space into subintervals.

¹⁰Also, an important technical difference is the cardinality of the state space. Because it is finite in Watson’s model, the set of all bijective mappings from reported states to ideal actions is also finite. In our setup, however, the state and action spaces are continuous. Applying Watson’s idea to this setup is equivalent to the uniform randomization over the space of all bijective functions from the interval of reported states to the interval of ideal actions, which is a large and mathematically complicated space.

buyer is privately informed about her type. Third, Krähmer’s and our constructions utilize conceptually different properties of private signal structures to extract the full information and surplus from the buyer. Signal structures in Krähmer (2020) are designed to monitor the buyer and detect his deviations from truth-telling. In particular, each signal structure is endowed with an individual signal set. This set is privately known to the seller, while the buyer privately observes the signal realization only. Thus, after receiving an ‘incorrect’ signal, the seller infers that the buyer lies and takes a penalizing action. A threat of this action enforces the buyer’s incentive-compatibility.¹¹ In our model, the signal structures share a common signal space, which implies that the buyer’s deviations are undetectable. However, a randomization between signal structures creates the uncertainty for the buyer about her payments contingent on the realized signal structure. As a result, signal distortions create a trade-off between her marginal benefits and losses. Fourth, a private signal structure in our model randomizes between two deterministic signal structures, whereas private signal structures in Krähmer (2021) are based on a randomization over a continuum of signal structures with individual signal spaces. Finally, our construction employs the continuity of the state space, whereas Krähmer’s approach relies on its discreteness.

Our paper is also related to the literature on the surplus extraction. This topic drew significant attention due to a seminal work by Crémer and McLean (1988), who demonstrated the possibility of full surplus extraction based on eliciting beliefs about buyers’ correlated values. On the other hand, recent developments in information design inspired by Kamenica and Gentzkow (2011) have demonstrated that it can also be a powerful tool for the seller to extract surplus from the buyer(s). Lewis and Sappington (1994), Johnson and Myatt (2006), Bergemann and Pesendorfer (2007), Esö and Szentes (2007), Li and Shi (2019), and Ivanov (2021) show that the seller(s) can benefit by designing or affecting buyers’ information about the product in standard bilateral trade or auctions.¹² Zhu (2021) and Larionov et al. (2021) consider the implementation problem with multiple (more than four) agents who can acquire additional information about their types. They demonstrate the possibility of implementing any social choice rule by using Shannon’s (1949) encryption technique and cross-checking agents’ messages. Pastrian (2021) demonstrates the full surplus extraction in the reduced form framework of McAfee and Reny (1992) with a behavioral subset of buyer’s types that are always truthful. Fu et al. (2021) consider a setup with a finite number of possible distributions of buyers’ values, where the seller has access to a finite number of independent draws from the true distribution. They establish that the full surplus extraction is feasible if the number of draws is large enough. Neither of these papers, however, considers private signal structures.

Finally, our construction of the private signal structure, which aligns the incentives of the better informed agent to reveal her information with those of the uninformed principal, exploits the idea of stochastically penalizing the agent for distortions in her behavior. Besides information design, a similar idea has been used in many other economic applications. For example, buyers’ truth-telling in the surplus-extracting mechanism by Crémer and McLean (1988) is sustained by payments that depend on the reports of all buyers. Because buyers’ valuations are correlated, distorting information by a buyer is more likely to induce profiles

¹¹Krähmer (2021) uses a similar idea in the cheap-talk context.

¹²Ivanov (2013) and Hwang et al. (2019) demonstrate the same result in competitive markets with horizontally differentiated products.

of reports that have low probabilities. Such profiles result in significant payments, which preclude the buyers' misreporting. Similarly, the literature on communication via a mediator utilizes the idea that a mediator can facilitate communication by stochastically penalizing the positively biased agent for reporting higher messages by recommending the principal to take lower actions (Goltsman et al., 2009; Ivanov, 2010). A similar effect can be achieved by adding the noise in the communication channel between the agent and the principal (Blume et al., 2007) or due to an imperfect ability of the principal to decode the agent's messages (Blume et al., 2019). Another application includes the principal-agent models with moral hazard. Ederer et al. (2018) study simple incentive schemes for an agent with privately known and asymmetric costs of performing different tasks. They show that optimal linear contracts include a randomization between payments rules with opposite reactions to performance on different tasks, since these contracts induce more balanced efforts and eliminate the efficiency losses from the agent's private information.

The rest of the paper is organized as follows. Section 2 introduces the general implementation model. Section 3 provides the main result for this framework. Section 4 applies these results to the bilateral-trade model. Finally, Section 5 concludes the paper.

2 Model

We consider the framework with two players, an agent (she) and a principal (he). The players communicate about the ex-ante unknown *state* $\theta \in \mathbb{R}$, which is a random variable drawn from the state space $\Theta = [\underline{\theta}, \bar{\theta}]$ according to a continuous density $f(\theta)$, such that $f > 0$. (Hereafter, $u > 0$ for a function $u : X \rightarrow Y$ means $u(x) > 0$ for all $x \in X$.) The principal's goal is to implement his *first-best* (or *ideal*) *action* $y(\theta)$ from a closed interval $\mathbf{A} \subset \mathbb{R}$. We assume that $y(\theta)$ is continuous and strictly increasing on $\Theta_0 = [\theta_0, \theta_1] \subset \Theta$, $y(\theta) = y_0$ for $\theta < \theta_0$, and $y(\theta) = y_1$ for $\theta > \theta_1$. Thus, the principal is interested in eliciting the perfect information about θ if it is in the subinterval of *target states* Θ_0 . For other states, knowing that $\theta < \theta_0$ or $\theta > \theta_1$ is sufficient to implement the principal's ideal actions y_0 or y_1 , respectively. As a normalization, it is without loss of generality to put $y(\theta) = \theta$ for $\theta \in \Theta_0$.¹³ By the continuity of $y(\theta)$, we have $y(\theta) = \theta_0$ for $\theta < \theta_0$ and $y(\theta) = \theta_1$ for $\theta > \theta_1$. Thus, Θ_0 is also a set of the principal's ideal actions induced by all states. We assume that $\Theta_0 \subset \mathbf{A}$, that is, all principal's ideal actions are feasible.

Two important comments about the principal's interests are worth mentioning. First, our model does not require specifying the principal's underlying preferences that are expressed, for example, by the payoff function $U(a, \theta)$, which has a maximizer $y(\theta)$ over a . For our purposes, knowing only $y(\theta)$ is sufficient. Thus, our results are robust to any modifications in the principal's payoff function as long as they do not affect the maximizer $y(\theta)$.¹⁴ The reason for this robustness, which is the second comment, is that we show that the principal can implement his ideal action $y(\theta)$ upon learning the state θ by the means of private information design. As a result, knowing the principal's preferences over alternatives (a, θ)

¹³Otherwise, if $y(\theta) \neq \theta$ for $\theta \in \Theta_0$, then the monotone transformation $z = y(\theta)$ results in $y(z) = z$.

¹⁴For example, consider $\mathbf{A} = \mathbb{R}_+$ and the principal's payoff functions $U_1(a, \theta) = -(a - \theta)^2$, $U_2(a, \theta) = \theta \ln(a) - a$, and $U_3(a, \theta) = 1$ if $a = \theta$ and 0 otherwise. Because $y(\theta) = \theta$ is the unique maximizer for all functions, the models with these principal's payoffs are equivalent.

different from the first-best outcomes $(y(\theta), \theta)$ is unnecessary.

The agent's payoff function $V(a, \theta, \gamma)$ depends on the principal's action a , the state θ , and the privately known *buyer's type* (or *type*) γ . In general, V can be independent of θ and/or γ . The type γ is drawn randomly from the type space \mathbf{T} according to the cdf $H(\gamma)$. We do not specify any assumptions about \mathbf{T} and $H(\gamma)$, however, there can be bounds on ordered \mathbf{T} for some specific V in order to guarantee that the payoff function respects the conditions imposed below. The variables θ and γ are independent.

We assume that $V(a, \theta, \gamma)$ is continuously differentiable in (a, θ) and $V'_a(a, \theta, \gamma) \neq 0$ for all $(a, \theta, \gamma) \in \Theta^2 \times \mathbf{T}$. The function $V'_a(a, \theta, \gamma)$ represents the agent's marginal payoff with respect to the principal's action. This function has another interpretation, which we employ through the paper. Because the principal wants to match the action to the state, $V'_a(a, \theta, \gamma)$ can be interpreted as the agent's marginal payoff from inducing the principal's (wrong) posterior belief that the state is a rather than θ . In other words, it represents the agent's marginal benefits or losses from manipulating the principal's action via his posterior belief concentrated at a . We impose the following separability condition on this function. Specifically, there exists a subinterval of states $\Theta_0 \subset \Theta$, such that $V'_a(a, \theta, \gamma)|_{a=\theta} = V'_a(\theta, \theta, \gamma)$ is given by

$$V'_a(\theta, \theta, \gamma) = g(\gamma) \zeta(\theta) \text{ for } (\theta, \gamma) \in \Theta_0 \times \mathbf{T}, \quad (1)$$

where $g(\gamma) > 0$ for all $\gamma \in \mathbf{T}$. That is, the agent's marginal payoff with respect to the action at the principal's ideal point $a = \theta$ can be factorized into functions $g(\gamma)$ and $\zeta(\theta)$, which solely depend on the agent's type γ and state θ , respectively, and do not change sign. Because $g > 0$, this implies that $V'_a(\theta, \theta, \gamma) \geq 0$ if and only if $\zeta(\theta) \geq 0$ for $\theta \in \Theta_0$.¹⁵ For concreteness, hereafter we assume that $\zeta(\theta) < 0$ for $\theta \in \Theta_0$. The case of $\zeta(\theta) > 0$ is symmetric.

Importantly, for a given state θ , the condition (1) is local as the factorization is required at the single point $a = \theta$ only. Intuitively, this condition requires the agent's marginal payoff be proportional to both the state and the agent's type at the principal's ideal action.

Information. A *signal structure* ξ determines a probability distribution $F_\xi(s|\theta)$ over signals s conditional on the state θ . For simplicity and with a minor abuse of notation, hereafter we use the term ξ for $F_\xi(s|\theta)$. A *signal set* $\mathbf{S}_\xi \subset \mathbb{R}$ is the support of ξ . A signal structure ξ is called a *signal function* if it maps each state $\theta \in \Theta$ into a signal $s = \xi(\theta)$. In this case, the signal set \mathbf{S}_ξ is the image of ξ . A signal function is *perfectly informative* if it is injective. Hereafter, we restrict the codomain \mathbf{C} of each signal function ξ by its image \mathbf{S}_ξ . Hence, a perfectly informative $\xi : \Theta \rightarrow \mathbf{C}$ is bijective and thus has the inverse function (hereafter called the *inverse*) $\varphi = \xi^{-1} : \mathbf{C} \rightarrow \Theta_\varphi$, where Θ_φ is the image of φ . Similarly to signal functions, we restrict the codomain of φ by its image Θ_φ . Because of the restrictions on the codomains of ξ and φ , the existence of a function $\xi : \Theta \rightarrow \mathbf{C}$ (or $\varphi : \mathbf{C} \rightarrow \Theta$) also implies that the image of ξ is \mathbf{C} (or Θ). Let \mathcal{I} be the space of all signal structures. A *private signal structure* $\rho \in \Delta\mathcal{I}$ is a probability distribution over signal structures whose realization ξ is privately observed by the principal. Denote $\rho(\xi)$ the probability of drawing ξ by ρ , and \mathcal{I}_ρ the support of ρ .

Timing. The game is played as follows. The agent is a priori uninformed about θ and

¹⁵The case of $\zeta(\theta) = 0$ is ruled out by $V'_a \neq 0$.

informed about γ . That is, her information about θ is determined by the prior density $f(\theta)$. At the beginning of the game, the principal publicly selects a private signal structure $\rho \in \Delta \mathcal{I}$ and an action $y_\xi(m)$.¹⁶ Then, the state θ and the signal structure $\xi \in \mathcal{I}_\rho$ are randomly and independently drawn according to f and ρ , respectively, where ξ becomes the private information of the principal.¹⁷ The agent then privately observes a signal s generated by ξ from θ and sends a message m from the message space \mathbf{M} to the principal who takes an action $y_\xi(m)$. Hereafter, we assume that $\mathbf{M} = \mathbf{S} = \bigcup_{\xi \in \mathcal{I}_\rho} \mathbf{S}_\xi$, that is, the message space is large enough to convey all information about signals.¹⁸

Conditional on ρ and $y_\xi(m)$, the following subgame is the decision problem with a privately and imperfectly informed agent. A strategy of the agent $m(s, \gamma, \rho) \in \Delta \mathbf{S}$ specifies a (possibly random) message m given her information: the observed signal $s \in \mathbf{S}$, the type γ , and the private signal structure ρ . An *optimal strategy* $m^*(s, \gamma, \rho)$ is a maximizer of the agent's posterior payoff

$$EV(m|s, \gamma, \rho) = \int_{\mathcal{I}_\rho} \int_{\Theta} V(y_\xi(m), \theta, \gamma) dq(\theta|s, \rho) d\rho(\xi), \quad (2)$$

where $q(\theta|s, \rho) \in \Delta \Theta$ is the agent's *posterior belief*, which is a probability distribution over θ derived from s and ρ by using Bayes' rule.¹⁹ We say that the state θ is *posterior and induced* by a signal s under a private signal structure ρ if θ is in the support of $q(\cdot|s, \rho)$. In particular, if the support \mathcal{I}_ρ of ρ contains only perfectly informative signal functions ξ , then $\theta = \varphi_\xi(s) = \xi^{-1}(s)$ represents the *posterior state* induced by a signal s under a signal function ξ , and thus the support of $q(\theta|s, \rho)$ is given by $\{\varphi_\xi(s) : \xi \in \mathcal{I}_\rho\}$.

Finally, the agent's *truthful* strategy is optimal under ρ if $m^*(s, \gamma, \rho) = s$ is in the set of maximizers of (2) for ρ and all $(s, \gamma) \in \mathbf{S} \times \mathbf{T}$. Equivalently,

$$EV(s|s, \gamma, \rho) = \max_{m \in \mathbf{S}} EV(m|s, \gamma, \rho) \text{ for all } (s, \gamma) \in \mathbf{S} \times \mathbf{T}. \quad (3)$$

Two comments are necessary here. First, since our model requires specifying the principal's ideal action $y(\theta)$ only, the nature of the principal's action $y_\xi(m)$ is unclear if he infers imperfect information about θ . However, it is not an issue as we focus on implementing $y(\theta)$ for all θ . Formally, we construct ρ , which randomizes between perfectly informative signals functions $\xi(\theta)$ only and sustains the truthtelling equilibrium. Hence, $y_\xi(s)|_{s=\xi(\theta)} = \theta$ for all $\xi \in \mathcal{I}_\rho$ and $\theta \in \Theta$. In other words, if a signal s is generated by any $\xi \in \mathcal{I}_\rho$ in state θ , and the agent reports s truthfully, then the principal infers θ and implements his ideal action $y(\theta) = \theta$. Second, the principal's ability to commit to $y_\xi(m)$ does not play a role. For

¹⁶Formally, the principal's $y_\xi(m)$ action can also be based on the private signal structure ρ . Because neither of our results is driven by this dependence, we omit it for simplicity of notation.

¹⁷Since the probability $\rho(\xi)$ does not depend on θ , ξ and θ are independent random variables. Hence, knowing ξ does not provide any additional information about θ to the principal.

¹⁸In general, the principal can specify the richer message space $\mathbf{M} = \mathbf{S} \times \mathbf{T}$ and request the agent to report information not only about the state, but also about her type. However, our main results do not rely on the principal's knowledge of the agent's type and thus do not require the agent report it.

¹⁹Since γ and θ are independent, $q(\theta|s, \rho)$ does not depend on γ .

example, he can commit to $y_\xi(m)$ ex-ante, or it can be interim-optimal, i.e., a solution to the problem of maximizing the principal’s expected payoff conditional on ξ and m .²⁰

3 Perfect information extraction and implementation

Before starting the general construction (hereafter, a *construction*) of private signal structures, which elicit the perfect information about the state from the agent and allow the principal to implement his ideal action, we provide an illustrative example. This example shows that perfect implementation can be achieved by randomizing between two simple and perfectly informative signal functions. Notably, we consider the agent’s preferences, which are independent of the state and the agent’s type. (The second example below considers general preferences, which depend on both the state and the type.) If the agent’s signal structure were publicly known, then the invariance of the agent’s preferences to information combined with the strict monotonicity of her payoff in the principal’s action would make the principal’s problem of eliciting any relevant information unsolvable. The insolvability of the problem is due to three factors. First, because the information about the state has no value to the agent, the principal cannot exploit the agent’s incentives to acquire this information.²¹ Second, the agent is more willing to share her information upon learning it if her preferences are closely aligned with those of the principal. Because the principal’s ideal action depends on the state, while the agent’s one does not, the players’ preferences are substantially conflicting. Finally, because the agent’s preferences are strictly monotone in action, then irrespectively of the information received, the agent will send only those messages that induce one of the extreme feasible actions.²² Together, these factors completely suppress the agent’s incentives to convey any relevant information. In contrast, assigning a private signal structure to the agent allows the principal to extract perfectly precise information from him.

3.1 Example A: action-only dependent agent’s preferences

Suppose that the state is uniformly distributed on the unit interval, i.e., $f(\theta) = 1, \theta \in \Theta = [0, 1]$ and the agent’s payoff function is of the form

$$V(a, \theta, \gamma) = V(a) = -a^b,$$

²⁰In the latter case, however, $y_\xi(m)$ is common knowledge in any equilibrium. That is, the agent expects the principal to play $y_\xi(m)$ in equilibrium.

²¹In the case of state-dependant agent’s preferences, the agent’s incentives to acquire information can play a critical role for information extraction. For example, Ivanov (2015, 2016) shows that the principal can elicit perfect information from the agent and implement her ideal actions in the cheap-talk framework by exploiting these incentives in a dynamic way.

²²Chakraborty and Harbaugh (2010) and Lipnowski and Ravid (2020) study cheap-talk communication with a perfectly informed agent whose preferences depend on the principal’s action only. As they show, the agent can disclose relevant information if the action set is multi-dimensional or the agent’s payoff function is not monotone in the action. In our example with single-dimensional actions and monotone preferences of the agent, informative communication is not feasible regardless of the agent’s public signal structure.

where $a \in \mathbf{A} = \mathbb{R}_+$ and $b \geq 1$ is a known parameter. Because V is strictly decreasing in a for $a \geq 0$, the agent's payoff is maximized at $a_0 = 0$.²³

Now, consider the private signal structure ρ^o , which randomizes with equal probabilities between two perfectly informative signal functions

$$\xi_1(\theta) = \theta \text{ and } \xi_2(\theta) = \left(1 - \theta^{\frac{b+1}{2}}\right)^{\frac{2}{b+1}}. \quad (4)$$

Because the images of functions ξ_1 and ξ_2 are identical and equal to $\mathbf{S} = [0, 1]$, then any agent's deviation from truthtelling is undetectable by the principal. On the other hand, the agent cannot infer the realized signal function and, thus, the state θ upon observing the signal s . In particular, a signal s generates the agent's posterior beliefs

$$q(\theta|s, \rho^o) = \Pr\{\theta|s, \rho^o\} = \begin{cases} \frac{1}{1+|\varphi_2'(s)|} & \text{if } \theta = \varphi_1(s), \\ \frac{|\varphi_2'(s)|}{1+|\varphi_2'(s)|} & \text{if } \theta = \varphi_2(s), \text{ and} \\ 0 & \text{if } \theta \notin \{\varphi_1(s), \varphi_2(s)\}, \end{cases} \quad (5)$$

where φ_i is the inverse of ξ_i :

$$\varphi_1(s) = \xi_1^{-1}(s) = s \text{ and } \varphi_2(s) = \xi_2^{-1}(s) = \left(1 - s^{\frac{b+1}{2}}\right)^{\frac{2}{b+1}}$$

Denote $q_i(s, \rho)$ the probability of the posterior state $\theta_i = \varphi_i(s)$:

$$q_i(s, \rho) = \Pr\{\theta = \varphi_i(s) | s, \rho\} = q(\varphi_i(s) | s, \rho) \quad (6)$$

If the principal believes that the agent is truthful, then a message $m \in \mathbf{S}$ induces the action

$$y_{\xi_i}(m) = \theta_i = \varphi_i(m), i = 1, 2$$

under the signal function ξ_i . Therefore, the agent posterior payoff (2) is given by

$$\begin{aligned} EV(m|s, \rho^o) &= q_1(s, \rho^o) V(y_{\xi_1}(m)) + q_2(s, \rho^o) V(y_{\xi_2}(m)) \\ &= q_1(s, \rho^o) V(\varphi_1(m)) + q_2(s, \rho^o) V(\varphi_2(m)) \\ &= -q_1(s, \rho^o) (\varphi_1(m))^b - q_2(s, \rho^o) (\varphi_2(m))^b. \end{aligned}$$

As an illustration, consider the quadratic payoff function $V(a) = -a^2$. Fig. 1 depicts the signal functions ξ_1, ξ_2 , and the posterior probability $q_1(s, \rho^o)$ for this function. In this case,

$$\begin{aligned} \varphi_2(s) = \xi_2(s) &= \left(1 - s^{3/2}\right)^{2/3}, \\ q_1(s, \rho^o) &= \frac{\left(1 - s^{3/2}\right)^{1/3}}{s^{1/2} + \left(1 - s^{3/2}\right)^{1/3}}, q_2(s, \rho^o) = \frac{s^{1/2}}{s^{1/2} + \left(1 - s^{3/2}\right)^{1/3}}, \text{ and} \end{aligned}$$

²³Formally, $V'_a(0) = 0$ for $b > 1$, which violates the condition $V'_a < 0$. However, since $V(a)$ is strictly decreasing in a , neither of our results is affected by this technicality.

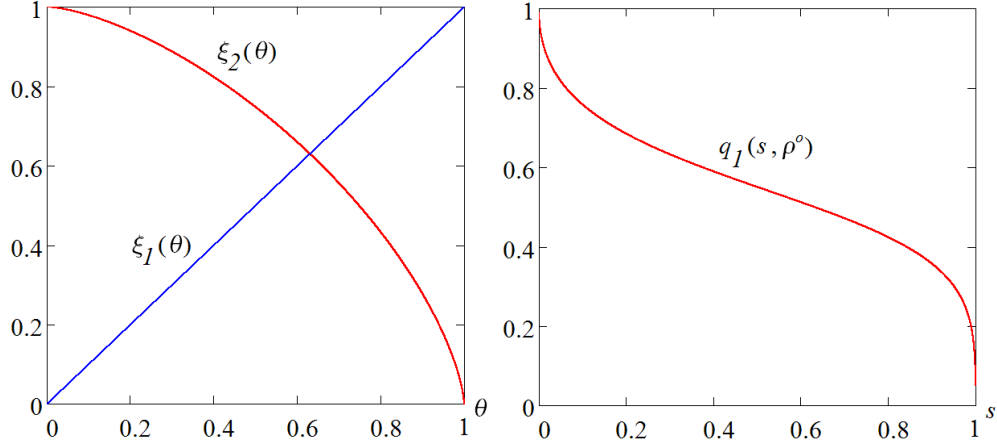


Figure 1: Signal functions $\xi_i(\theta)$, $i = 1, 2$ and the posterior probability $q_1(s, \rho^o)$ for $f(\theta) = 1$ and $V(a) = -a^2$.

$$EV(m|s, \rho^o) = -\frac{(1 - s^{3/2})^{1/3} m^2 + s^{1/2} (1 - m^{3/2})^{4/3}}{s^{1/2} + (1 - s^{3/2})^{1/3}}.$$

By using simple calculus, it is easy to verify that $EV(m|s, \rho^o)$ is maximized at unique $m = s$ for all $s \in \mathbf{S}$, i.e., the agent strictly benefits from reporting truthfully.

Intuitively, this example demonstrates the keys factors of private signal structures, which sustain agent's truthtelling. Specifically, the induced posterior states $\theta_1 = \varphi_1(m)$ and $\theta_2 = \varphi_2(m)$ associated with signal structures ξ_1 and ξ_2 , respectively, react oppositely to an agent's message m . Next, the principal's ideal action $y(\theta_i) = \theta_i$ is monotone in the induced posterior state θ_i , $i = 1, 2$. Finally, the agent's payoff $V(a)$ is monotone in the principal's action a . Together, these factors create the trade-off for the agent: any distortion of the signal in an attempt to marginally benefit from the receiver's action taken under one signal function are offset by the marginal losses caused by the action taken under the other signal function.

At the same time, the magnitude of the trade-off between the agent's marginal benefits and losses from distortions is driven by the shapes of signal functions, specifically, their inverses φ_1 and φ_2 . Their effect on the trade-off is dual. First, they determine the marginal effects of an agent's message on the agent's posterior payoff via actions taken under different signal structures. Second, they reallocate the agent's posterior beliefs between the posterior states. We now explain in detail the relationship between the shapes of the inverses and their overall effect on the agent's incentives to report truthfully.

In order to explain the first effect, recall that the principal's action $y_\xi(m)$ taken for a signal function ξ matches the induced posterior state $\theta_i = \varphi_i(m)$. As a result, the shape of $\varphi_i(m)$ determines the marginal effect of an agent's message m on the principal's action a_i . Next, note that the agent's payoff function $V(a)$ is strictly concave in the principal's action a for $b > 1$. Hence, the agent's marginal benefits $|V'(a)|$ from inducing a lower action are larger for high actions. As a result, the overall effect of signal distortions on the agent's posterior payoff depends on the interaction between the inverses φ_1 and φ_2 and the marginal payoff $V'(a)$. Specifically, note that the first inverse φ_1 is linear in s . Therefore, the marginal

effect of the signal s on the principal's action $a_1 = \varphi_1(s)$ is constant, $\varphi_1'(s) = 1$. On the other hand, because φ_2 is strictly concave, the absolute value of the marginal effect $|\varphi_2'(s)|$ on $a_2 = \varphi_2(s)$ is increasing in s . This implies that the ‘counter-moving’ action a_2 increases at the faster rate in response to downward distortions if the signal s is high, whereas the rate of a decrease in the ‘co-moving’ action a_1 is constant. In other words, the marginal penalty from downward distortions caused by the penalizing action a_2 relative to the benefits from the favorable action a_1 is increasing in an agent's signal s . Because the agent has the stronger incentive to decrease the action a_1 if it is high and, thus, is associated with a high signal s , then understating such signals will result in the higher losses from the penalizing action a_2 .

Second, the shapes of the inverses reallocate the agent's posterior beliefs between the posterior states $\theta_1 = \varphi_1(s)$ and $\theta_2 = \varphi_2(s)$ from states in which the agent has the stronger incentives to lie toward those with weaker incentives. In our example, the buyer's payoff is steeper for high actions. Because the principal matches actions to states, the agent's benefits from distorting information are higher for high states. As a result, for high signals the slopes $|\varphi_1'|$ and $|\varphi_2'|$ of the inverses φ_1 and φ_2 assign a lower posterior probability to the higher posterior states. Because the agent places lower probabilities on these states, this decreases her posterior (that is, conditional on signal) benefits from manipulating the signal and thus her incentives to lie. The balance between these two forces sustains agent's truthtelling.

Given these observations, it is easy to notice the complementarity between the marginal effects of actions and their probabilities on the agent's incentives to report truthfully. A steeper slope $|\varphi_2'(s)|$ in response to signal s results in both the higher marginal penalty from the counter-moving action a_2 and the higher probability $q_2(s, \rho)$ of inducing this action. In other words, the higher magnitude of one effect intensifies the second effect as well. Next, it is worth noting that the agent remains truthful regardless of the precision of her information (measured, for instance, by the variance of posterior states). In fact, if $\hat{s} = 2^{-2/3} \simeq 0.63$, then the agent is perfectly informed about the posterior state $\hat{\theta} = \hat{s}$. However, she still cannot use this information in her favor.

Finally, note that perfect information extraction does not rely on the strict concavity of the agent's preferences. Specifically, the main arguments above hold through if the agent's payoff is linear in a , that is, $V(a) = -a$. However, there are two key differences between the cases of the linear and strictly concave payoff functions. First, the agent with the linear payoff function is indifferent among all messages m for any signal $s \in \mathbf{S}$. At the same time, the example above demonstrates that the agent's incentive to report truthfully are not driven by this indifference. Second, (4) implies that the inverse φ_2 is linear for $V(a) = -a$. As a result, the posterior beliefs $q_i(s, \rho^o) = \frac{1}{2}, i = 1, 2$ given by (5) are constant. However, the linearity of the inverses does not allow us to see an interplay between their shapes, the posterior beliefs, and the agent's marginal benefits and losses from distorting her information.

3.2 Optimal private signal structures

We start the general construction with the following lemma. It demonstrates how the agent's posterior beliefs are shaped by a private signal structure, which randomizes between two perfectly informative signal functions.

Lemma 1 (*Ivanov and Sam, 2022*) *Consider a private signal structure ρ , which randomizes between differentiable signal functions $\xi_1 : \Theta \rightarrow \mathbf{S}$ and $\xi_2 : \Theta \rightarrow \mathbf{S}$ with probabilities $p_1 \in$*

$(0, 1)$ and $p_2 = 1 - p_1$, respectively, where $\mathbf{S} = [s_0, s_1]$, $s_1 > s_0$, and $\xi_i' \neq 0$. Denote $\varphi_i = \xi_i^{-1}$ the inverse of ξ_i . Then,

$$q_i(s, \rho) = \frac{p_i f(\varphi_i(s)) |\varphi_i'(s)|}{p_1 f(\varphi_1(s)) |\varphi_1'(s)| + p_2 f(\varphi_2(s)) |\varphi_2'(s)|}. \quad (7)$$

Intuitively, the lemma highlights the key feature of signal functions, specifically, the possibility to induce the agent's posterior beliefs $q_i(s, \rho)$ about posterior states $\theta_1 = \varphi_1(s)$ and $\theta_2 = \varphi_2(s)$ anywhere between 0 and 1 by varying the ratio $\frac{|\varphi_2'(s)|}{|\varphi_1'(s)|}$ of the slopes of the inverses φ_i' . To see this feature, suppose $p_1 = p_2 = \frac{1}{2}$ and f is uniform, i.e., $f(\theta) = \frac{1}{\theta - \bar{\theta}}$. It follows then that $q_1(s, \rho) = \frac{1}{1 + \frac{|\varphi_2'(s)|}{|\varphi_1'(s)|}}$ and $q_2(s, \rho) = 1 - q_1(s, \rho)$. By varying the ratio

$\frac{|\varphi_2'(s)|}{|\varphi_1'(s)|}$ between 0 and ∞ , the principal can induce $q_i, i = 1, 2$ anywhere between 0 and 1. As shown in the example above, the principal can use this ratio in order to reallocate the agent's posterior beliefs from the states with higher benefits from distorting the signal to the states with lower benefits and, as a result, reduce the agent's interim incentives to lie.

General construction. Consider the signal space $\mathbf{S} = [\underline{s}, \bar{s}]$ and a private signal structure ρ , which randomizes with equal probabilities between two perfectly informative signal functions ξ_1 and ξ_2 with the inverses $\varphi_i = \xi_i^{-1} : \mathbf{S} \rightarrow \Theta$ defined as follows. First, select a differentiable $\varphi_1 : \mathbf{S} \rightarrow \Theta$, such that $\varphi_1' > 0$. Thus, $\varphi_1(\underline{s}) = \underline{\theta}$ and $\varphi_1(\bar{s}) = \bar{\theta}$. The principal's problem is to derive $\varphi_2 : \mathbf{S} \rightarrow \Theta$, such that the private signal structure ρ sustains the agent's truth-telling and, thus, allows the principal to implement $y(\theta)$ upon inferring θ from m and ξ_i . Importantly, the signal sets, that is, the images of functions ξ_1 and ξ_2 are identical and equal to \mathbf{S} . First, this implies that the agent is unable to infer the realized ξ_i upon observing the signal s . Second, any agent's deviation from truth telling is undetectable by the principal.

Given the agent's truthful strategy $m^*(s, \gamma, \rho) = s$ for $(s, \gamma) \in \mathbf{S} \times \mathbf{T}$, the principal's best response to message m under the signal structure ξ_i is

$$y_{\xi_i}(m) = \varphi_i(m), i = 1, 2. \quad (8)$$

The agent's problem upon receiving a signal s is to maximize her posterior payoff (2) over messages $m \in \mathbf{S}$. Using (8), the posterior payoff can be expressed as

$$EV(m|s, \gamma, \rho) = \sum_{i=1}^2 q_i(s, \rho) V(\varphi_i(m), \varphi_i(s), \gamma) \text{ for } (m, s, \gamma) \in \mathbf{S}^2 \times \mathbf{T}. \quad (9)$$

Then, the agent's marginal posterior payoff is given by

$$\frac{\partial}{\partial m} EV(m|s, \gamma, \rho) = \sum_{i=1}^2 q_i(s, \rho) V'_a(\varphi_i(m), \varphi_i(s), \gamma) \varphi_i'(m), \quad (10)$$

and truthtelling is optimal if and only if

$$EV(s|s, \gamma, \rho) = \max_{m \in \mathbf{S}} EV(m|s, \gamma, \rho) \text{ for all } (s, \gamma) \in \mathbf{S} \times \mathbf{T}. \quad (11)$$

By using (10), the first-order condition for the agent's maximization problem (11) is

$$\frac{\partial}{\partial m} EV(m|s, \gamma, \rho)|_{m=s} = \sum_{i=1}^2 q_i(s, \rho) V'_a(\varphi_i(s), \varphi_i(s), \gamma) \varphi'_i(s) = 0, (s, \gamma) \in \mathbf{S} \times \mathbf{T}. \quad (12)$$

Next, invoking condition (1) results in

$$V'_a(\varphi_i(s), \varphi_i(s), \gamma) = g(\gamma) \zeta(\varphi_i(s)),$$

where $g > 0$ and $\zeta > 0$. This means that (12) is independent of γ and thus can be written as

$$\frac{\partial}{\partial m} EV(m|s, \rho)|_{m=s} = \sum_{i=1}^2 q_i(s, \rho) \zeta(\varphi_i(s)) \varphi'_i(s) = 0 \text{ for all } s \in \mathbf{S}.$$

By using Lemma 1, (12) can be written as a separable differential equation with respect to φ_2 for a given φ_1 :

$$\varphi'_1(s) |\varphi'_1(s)| f(\varphi_1(s)) \zeta(\varphi_1(s)) + \varphi'_2(s) |\varphi'_2(s)| f(\varphi_2(s)) \zeta(\varphi_2(s)) = 0, \quad (13)$$

Because $\varphi'_1 = |\varphi'_1| > 0$, $f > 0$, and $\zeta > 0$, it immediately follows that $\varphi'_2 < 0$ and, hence, $|\varphi'_2(s)| = -\varphi'_2(s)$. Therefore, (13) can be expressed as

$$\varphi'_1(s) h(\varphi_1(s)) = -\varphi'_2(s) h(\varphi_2(s)), \quad (14)$$

where

$$h(\theta) = \sqrt{f(\theta) |\zeta(\theta)|} = \sqrt{-f(\theta) \zeta(\theta)} > 0.$$

The solution to (14) with the boundary condition $\varphi_2(\underline{s}) = \bar{\theta}$ is

$$\varphi_2(s) = \Psi^{-1}(\Psi(\underline{\theta}) + \Psi(\bar{\theta}) - \Psi(\varphi_1(s))), \quad (15)$$

where $\Psi(x) = \int h(x) dx$ is the antiderivative of h .²⁴

In general, a pair of inverses $\varphi_i : \mathbf{S} \rightarrow \Theta, i = 1, 2$ related by (15) is not necessarily a solution to the agent's maximization problem (11) as the second-order conditions might not hold. The following regularity condition addresses this issue.

Condition 1 Given $\Theta_0 \subset \Theta$, $\nu(a, \theta, \gamma) = \frac{V'_a(a, \theta, \gamma)}{h(a)}$ is decreasing in a for all $(a, \theta, \gamma) \in \Theta_0^2 \times T$.

Notably, this condition is imposed on the model primitives only: the payoff function $V(a, \theta, \gamma)$ and the prior density $f(\theta)$. Therefore, agent's truthtelling can be optimal for various pairs $\{\varphi_1, \varphi_2\}$ parameterized by the inverse φ_1 . Condition 1 can be explained by

²⁴Note that $\varphi_1(\bar{s}) = \bar{\theta}$ implies $\varphi_2(\bar{s}) = \underline{\theta}$, which means that functions φ_1 and φ_2 have identical images \mathbf{S} .

noting that $\nu(a, \theta, \gamma)$ is the ratio of two functions, $V'_a(a, \theta, \gamma)$ and $h(a) = \sqrt{f(a)|\zeta(a)|}$. The first function $V'_a(a, \theta, \gamma)$ is the marginal payoff with respect to the principal's action or equivalently, the induced posterior state. It is decreasing in a if $V(a, \theta, \gamma)$ is concave in a , and increasing if $V(a, \theta, \gamma)$ is convex in a . The second function $h(a)$ reflects the marginal benefit $|\zeta(a)|$ from distorting an action a weighted by the prior density $f(a)$. In this light, Condition 1 requires the function $V(a, \theta, \gamma)$ be 'not very convex' in a , and the agent's weighted marginal benefit $h(a)$ be 'relatively decreasing' in action a (since $V'_a < 0$ and $h > 0$ imply $\nu < 0$).²⁵

Given these preliminaries, the following theorem establishes the main result of the paper. Consider the private signal structure ρ^* , which randomizes with equal probabilities between signal functions $\xi_1 = \varphi_1^{-1}$ and $\xi_2 = \varphi_2^{-1}$, such that the relationship between φ_1 and φ_2 is given by (15). Then ρ^* allows the principal to elicit the perfect information about the state from the agent and, hence, implement his ideal action if the above regularity condition holds.

Theorem 1 *Suppose V satisfies (1) and (f, V) satisfy Condition 1 for $\Theta_0 = \Theta$. Consider a pair of functions (φ_1, φ_2) related by (15), where $\varphi_1 : \mathbf{S} \rightarrow \Theta$ is differentiable and $\varphi'_1 > 0$. Then the private signal structure ρ^* that randomizes between φ_1^{-1} and φ_2^{-1} with equal probabilities sustains agent's truthtelling for all $(\theta, \gamma) \in \Theta \times \mathbf{T}$.*

This result extends Theorem 1 in Ivanov and Sam (2022) in three dimensions. First, in their model the agent is ex-ante uninformed, while our setup allows the agent to be privately informed about her type γ . Second, their model assumes the existence of the agent's ideal action $y_A(\theta)$ for each state θ . Third, they assume the strict supermodularity of the agent's preferences in (a, θ) . This implies the dependence of the agent's payoff on state θ . Our setup does not require the existence of the agent's ideal action, the dependence of her payoff function on the state, or the supermodularity. All these assumptions are replaced with the strict monotonicity of the agent's payoff in the principal's action. Specifically, recall that the opposite monotonicities of inverses $\varphi_i(s), i = 1, 2$ in signal s and the monotonicity of the principal's action $y(\theta) = \theta$ in θ imply that the principal's actions $a_i = \varphi_i(m), i = 1, 2$ under different signal functions ξ_i react oppositely in response to the agent's message m . Then, the monotonicity of the agent's payoff V in the principal's action implies that the opposite reactions of $\varphi_i(m), i = 1, 2$ to m are mapped in the opposite marginal payoffs to the agent. That is, agent's misreporting in an attempt to obtain extra gains under one signal function are offset by the extra losses under the other signal function. These marginal effects are balanced by the relationship (15) between the inverses $\varphi_i(s), i = 1, 2$ in order to sustain agent's truthtelling. Notably, the logic above does not depend on the concavity of the agent's preferences in a . As a result, the theorem may hold for arbitrarily convex payoff functions V as long as (V, f) satisfy Condition 1.²⁶

²⁵If $V'_a > 0$, then $\nu > 0$. In this case, Condition 1 implies that $h(a)$ must be 'relatively increasing' in a .

²⁶Consider $V(a) = e^{-ba}$, which is decreasing in a , and the parameter b is the Arrow-Pratt measure of absolute risk aversion $R(a) = -\frac{V''_a}{V'_a} = b$. According to this measure, V becomes more convex as b increases and eventually converges to the extreme case of convexity: $V(0) = 0$ and $V(\theta) = -1$ for $\theta > 0$. Then $V'_a = -be^{-ab}$, and $\zeta(\theta) = -V'_a(a)|_{a=\theta} = be^{-b\theta} > 0$. Also, consider the truncated exponential distribution on $[0, 1]$ with the density $f(\theta) = c_0 e^{c\theta}$, where $c_0 = \frac{c}{1-e^{-c}}$. Then $h(a) = \sqrt{\zeta(a)f(a)} = \sqrt{c_0 b e^{-\frac{b-c}{2}a}}$, and $\nu(a) = \frac{V'_a(a, \theta, \gamma)}{h(a)} = -\sqrt{\frac{b}{c_0}} e^{\frac{c-b}{2}a}$ is decreasing in a for any b and $c > b$.

At the same time, the proofs of the two theorems share common features. In both theorems, the critical part is to establish the optimality of the agent's truth-telling strategy under the private signal structure ρ . Once the agent's posterior payoff function (9) is single-peaked in m and achieves its maximum at $m = s$, this prevents the agent's local distortions of her signal (i.e., small lies) as well as global distortions (large lies). In order to guarantee the single-peakedness of the posterior payoff, we employ the pseudo-concavity property.²⁷ The tension with pseudo-concavity, however, comes from three factors. First, while each function $V(\varphi_i(m), \varphi_i(s), \gamma), i = 1, 2$ in (9) is strictly monotone and, hence, pseudo-concave in m , the agent's posterior payoff EV is a convex combination of these functions, which is generally not pseudo-concave.²⁸ Second, a function $V(\varphi_i(m), \varphi_i(s), \gamma), i = 1, 2$ is a composite function of V and φ_i . As a result, the pseudo-concavity of this function is violated if the (local) concavity of V in a is dominated by the convexity of $\varphi_i(m)$. Third, φ_1 and φ_2 are functionally dependent by (15). Therefore, the posterior payoff EV is pseudo-concave in m if: i) each composite function $V(\varphi_i(m), \varphi_i(s), \gamma), i = 1, 2$ is pseudo-concave, and ii) a convex combination of these functions is also pseudo-concave.

Condition 1 resolves all three issues. Specifically, the necessary and sufficient condition for the pseudo-concavity of $EV(m|s, \gamma, \rho)$ is the pseudo-monotonicity of the marginal posterior payoff $\frac{\partial}{\partial m}EV(m|s, \gamma, \rho)$ (Hadjisavvas et al., 2005).²⁹ To show that this function is pseudo-monotone, we use the results by Quah and Strulovici (2012) who establish conditions for the pseudo-monotonicity of a convex combination of pseudo-monotone functions. We apply these conditions to composite functions $V(\varphi_i(m), \varphi_i(s), \gamma), i = 1, 2$, and use the functional relationship (14) between φ_1 and φ_2 . This completes the proof of the theorem.

4 Application: bilateral trade

As the leading economic application of the results above, we consider the bilateral trade model with generalized players' preferences. Specifically, the buyer's preferences are non-quasilinear and generally depend on two random variables, the state and the buyer's type. The buyer is privately informed about her type γ , which affects her payoff only. The seller determines the buyer's information about the state θ , which reflects the product quality and, thus, buyer's willingness to pay for the object. As the main result, we demonstrate that the seller can extract the perfect information about the state and, as a consequence, the full buyer's surplus by using a private signal structure that randomizes between two perfectly informative signal functions.

²⁷The pseudo-concavity is a generalized form of the concavity of differentiable functions for which stationary points are also global maximizers.

²⁸For example, $2e^m$ and $2e^{-m}$ are strictly monotone (i.e., pseudo-concave) in m , but $e^m + e^{-m}$ is not pseudo-concave.

²⁹A function $\phi(x)$ is *pseudo-monotone* on a convex set $\mathbf{A} \subset \mathbb{R}$ if for every $(x, y) \in \mathbf{A}^2$, $\phi(x)(y - x) \leq 0$ implies $\phi(y)(y - x) \leq 0$. Equivalently, $\phi(x) \leq 0$ implies $\phi(y) \leq 0$ for all $y > x$.

4.1 Setup

A buyer (she) and a seller (he) are involved in trading a single indivisible object. The buyer’s utility from obtaining the object and making a payment t to the seller is determined by the payoff function $V(t, \theta, \gamma)$. The state θ reflects the intrinsic characteristics of the object, for example, its quality. In our setup, it also entirely determines the buyer’s willingness to pay. That is, $t = \theta$ is the highest price that the buyer is willing to pay in state θ . In general, the state θ can also affect the seller’s payoff.³⁰ Similarly to the main framework, γ is the buyer’s type, which affects the buyer’s payoff only. This variable has some antecedents in the mechanism design literature. For instance, in Dworzak et al. (2021), a privately known variable reflects the marginal value for money of agents in a market. In our model, the meaning of γ is broader. As shown below, it can reflect the marginal (dis-)utility with respect to the state θ and payment t similarly to Dworzak et al. (2021). In addition, it can determine the concavity of the buyer’s payoff function (i.e., the magnitude of the risk-aversion) or other characteristics.

The state θ is a random variable drawn according to a continuous density $f(\theta) > 0$ from the state space $\Theta = [0, \bar{\theta}]$. The type γ is drawn randomly from the type space \mathbf{T} according to the cdf $H(\gamma)$. The variables θ and γ are independent. The statistical independence between the ‘quality’ characteristics of the object and some intrinsic buyer’s preferences also has antecedents in the mechanism design literature, for example, Esö and Szentes (2007). According to them: “for a specific example, suppose that the object for sale is a car, and assume that the buyer knows its make, model, age, and mileage, but not its colour, which the seller can reveal. It seems reasonable to assume that a buyer’s initial willingness to pay for the car and his colour preference are statistically independent”.

Seller’s preferences. As in the main framework, we do not define the seller’s utility function $U(t, \theta)$. Instead, we assume that the seller’s goal is to extract the full surplus θ from the buyer in *target* states $\Theta_0 = [\theta_0, \bar{\theta}]$, where $\theta_0 \in \Theta$. Intuitively, Θ_0 is the subset of states in which the buyer’s willingness to pay, which is represented by θ , exceeds the seller’s utility from keeping the object.³¹

Buyer’s preferences. The buyer’s payoff function $V(t, \theta, \gamma)$ is continuously differentiable in (t, θ) and $V'_t(t, \theta, \gamma) < 0 < V'_\theta(t, \theta, \gamma)$ for all $(t, \theta, \gamma) \in \Theta^2 \times \mathbf{T}$. As a normalization, we assume that the value of the buyer’s outside option is 0, which she receives in the case of not obtaining the object and not making a payment. As noted above, θ represents the buyer’s willingness to pay, i.e., the maximum payment, which makes her indifferent between making this payment for the object and taking the outside option. That is, $V(t, \theta, \gamma)$ must satisfy the condition

$$V(\theta, \theta, \gamma) \equiv 0 \text{ for all } (\theta, \gamma) \in \Theta \times \mathbf{T}. \quad (16)$$

Furthermore, for a given $\theta \in \Theta$, $V'_t(t, \theta, \gamma) < 0$ implies that (16) holds only for $t = \theta$.

³⁰For example, if there is another market in which buyers’ valuations depend on the product quality, then this quality determines the seller’s expected opportunity cost and, hence, his payoff from keeping the object.

³¹As a very special case, suppose that the seller’s preferences are quasilinear, $U(t, \theta) = u(\theta) + t$, where $u(\theta)$ is the seller’s payoff from keeping the object in state θ , such that $u(\theta) \leq \theta$ if and only if $\theta \geq \theta_0$ for $\theta_0 \in \Theta$. As a result, the seller’s goal is to extract full surplus if and only if $\theta \in \Theta_0 = [\theta_0, \bar{\theta}]$.

Condition (16) implies that the buyer's willingness to pay does not depend on γ . This condition is not novel and imposed, for instance, in the literature on auctions with buyers' non-linear preferences (see Section 4.1 in Krishna, 2009). Specifically, this literature considers buyers with payoff functions $V(t, \theta) = v(\theta - t)$, where $v(\cdot)$ is in the class of functions $\mathbf{V} = \{\mathbf{C}^1 | v(0) = 0, v' > 0\}$.³² Thus, if V is an element of a subset $\mathbf{T} \subset \mathbf{V}$ parameterized by variable γ ,

$$V(t, \theta, \gamma) = v(\theta - t, \gamma), \quad (17)$$

then it must satisfy (16). For such payoff functions, while γ does not affect the maximum acceptable payment $t = \theta$ of the buyer, it affects her payoff for $t \neq \theta$. For example, the buyer's willingness to pay for a used car can be determined entirely by its quality characteristics such as the model, mileage, and age. At the same time, the marginal pleasure of driving the car is also affected by its color or the shape of her payoff function.

Finally, we assume that separability condition (1) holds for the subset of target states Θ_0 . That is, the buyer's marginal payoff with respect to the payment at point $t = \theta$ can be expressed as

$$V'_t(t, \theta, \gamma) |_{t=\theta} = g(\gamma) \zeta(\theta) \text{ for } (\theta, \gamma) \in \Theta_0 \times \mathbf{T}, \quad (18)$$

where $g(\gamma) > 0$ for all $\gamma \in \mathbf{T}$. Because $V'_t < 0$, it follows that $\zeta(\theta) < 0$ for all $\theta \in \Theta_0$.

Trade. The terms of trade are enforced by a trading mechanism (hereafter, a *mechanism*) \mathcal{M} defined as follows. A mechanism $\mathcal{M} = (\mathbf{M}, Q_\xi(m), t_\xi(m))$ consists of a message space \mathbf{M} , an allocation rule $Q_\xi(m) \in [0, 1]$, and a transfer rule $t_\xi(m) \geq 0$. Here, Q and t are the buyer's probability of obtaining the object and her payment to the seller, respectively. Importantly, Q and t depend on both the buyer's message m and the realized structure ξ privately known to the mechanism.³³

Timing. The game is played as follows. The buyer is a priori perfectly informed about γ and uninformed about θ . At the beginning of the game, the seller publicly selects a private signal structure $\rho \in \Delta \mathcal{I}$ and a mechanism $\mathcal{M} = (\mathbf{M}, Q_\xi(m), t_\xi(m))$. Then, the state θ and the signal structure $\xi \in \mathcal{I}_\rho$ are randomly drawn according to f and ρ , respectively, where ξ becomes the private information of the mechanism.³⁴ The buyer then privately observes a signal s generated by ξ from θ and decides whether to participate in trade or take the outside option. In the former case, the buyer sends a message $m \in \mathbf{M}$ to the mechanism \mathcal{M} . Finally, the terms of trade are enforced by the mechanism.

Because ρ is publicly observable, the following subgame is a standard selling mechanism with a privately and imperfectly informed buyer. Thus, we can invoke the Revelation Principle and restrict attention to direct interim incentive-compatible mechanisms, that is, such that $\mathbf{M} = \mathbf{S} = \bigcup_{\xi \in \mathcal{I}_\rho} \mathbf{S}_\xi$ and the buyer is truthful for all signals generated by the private signal structure ρ . Direct mechanisms are denoted $\mathcal{M} = (Q_\xi(s), t_\xi(s))$ hereafter. We also require mechanisms be interim individually-rational. This implies that the buyer does not

³²For risk-averse buyers, the concavity condition is additionally imposed $v'' < 0$.

³³In general, Q and t may also depend on the private signal structure ρ . Since our results do not rely on this dependence, we omit it in notation.

³⁴Since the probability $\rho(\xi)$ does not depend on θ , ξ and θ are independent random variables. Hence, knowing ξ does not provide any additional information about θ to the seller.

receive a negative posterior payoff from trade upon receiving any signal (since the value of her outside option is normalized to 0). Thus, hereafter we consider only those mechanisms, which are interim incentive-compatible (IC) and interim individually-rational (IR), i.e., those which satisfy the following conditions:

$$EV_B(s|s, \gamma, \rho) = \max_{m \in \mathbf{S}} EV_B(m|s, \gamma, \rho) \text{ for all } (s, \gamma) \in \mathbf{S} \times \mathbf{T}, \text{ (IC)} \quad (19)$$

$$EV_B(s|s, \gamma, \rho) \geq 0 \text{ for all } (s, \gamma) \in \mathbf{S} \times \mathbf{T}, \text{ (IR)} \quad (20)$$

Here, $EV_B(m|s, \gamma, \rho)$ is the buyer's posterior payoff

$$EV_B(m|s, \gamma, \rho) = \int_{\mathcal{I}_\rho} \int_{\Theta} Q_\xi(m) V(t_\xi(m), \theta, \gamma) + (1 - Q_\xi(m)) V(t_\xi(m), 0, \gamma) dq(\theta|s, \rho) d\rho(\xi),$$

where $V(t, 0, \gamma)$ is the buyer's value in the case of paying t for the worthless object, i.e., the one for which the buyer's willingness to pay is $\theta = 0$. It is equivalent to not obtaining the object while still paying t .

Given this framework, we start with an example, which provides the key insights into the general construction of private signal structures and mechanisms that extract full information and surplus from the buyer.

4.2 Example B: non-quasilinear players' preferences

Suppose that the prior density of states is uniform on the unit interval, that is, $f(\theta) = 1, \theta \in \Theta = [0, 1]$. The type $\gamma \in \mathbf{T}$ is drawn randomly according to the cdf $H(\gamma)$. The variables θ and γ are independent. The subset of seller's target states is $\Theta_0 = [\theta_0, 1]$.

Following the literature on auctions with risk-averse buyers, the buyer's payoff from consuming the product and making a payment t is given by the function (17):

$$V(t, \theta, \gamma) = v(\theta - t, \gamma),$$

where $v(x, \gamma)$ is strictly increasing, continuously differentiable, and concave in x , and $v(0, \gamma) = 0$ for all $\gamma \in \mathbf{T}$. This property and (17) imply that the buyer's willingness to pay is θ for all $(\theta, \gamma) \in \Theta \times \mathbf{T}$. That is, (17) satisfies condition (16). Furthermore, (17) implies

$$V'_t(t, \theta, \gamma)|_{t=\theta} = -v'_x(0, \gamma),$$

that is, condition (18) holds, where $g(\gamma) = v'_x(0, \gamma) > 0$ and $\zeta(\theta) = -1$.

Next, we verify that Condition 1 also holds. First, note that $h(\theta) = \sqrt{-f(\theta)\zeta(\theta)} = 1$. Second, because $v(x, \gamma)$ is concave in x , then $V'_t(t, \theta, \gamma) = -v'_x(\theta - t, \gamma)$ is decreasing in t . As a result, the function $\nu(t, \theta, \gamma) = -\frac{v'_x(\theta - t, \gamma)}{h(t)}$ is decreasing in t as well.

An example of function (17) is the linear-quadratic function

$$v(x, \gamma) = \gamma x - x^2,$$

where $\gamma > 2$. In this case, the buyer's type γ determines the relative weight of the linear component in the payoff and, hence, the marginal payoff with respect to the difference $\theta - t$

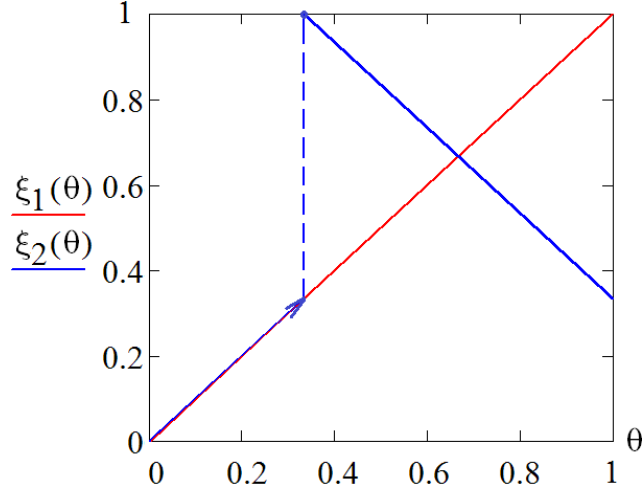


Figure 2: Signal functions $\xi_i(\theta)$, $i = 1, 2$ for $\theta_0 = \frac{1}{3}$.

at $t = \theta$. Another example is the hyperbolic absolute risk aversion (HARA) payoff function

$$v(x, \gamma) = \begin{cases} \frac{1-\gamma}{\gamma} \left(\frac{\alpha}{1-\gamma}x + \beta \right)^\gamma - \frac{1-\gamma}{\gamma} \beta^\gamma & \text{if } \gamma \neq 0, \\ \ln \left(1 + \frac{\alpha}{\beta}x \right) & \text{if } \gamma = 0, \end{cases}$$

where $\alpha > 0$ and $\frac{\alpha}{1-\gamma}x + \beta > 0$.³⁵ Depending on values of γ , α , and β , this form encompasses many standard payoff functions, such as linear, exponential (constant absolute risk aversion), power (constant relative risk aversion), and logarithmic.³⁶ For this payoff function, the value of γ determines the degree of buyer's risk aversion and thus substantially affect her utility and incentives in different mechanisms without the full surplus extraction. Because γ is the buyer's private information, the seller can be uncertain about the specific shape of the payoff function in our model. For example, he might not know whether it takes the power, exponential, or linear form. As we demonstrate below, the structure of the private signal structure and the mechanism are invariant to this information. At the same time, these tools jointly allow the seller to extract perfect information and full surplus from the buyer.

Now, consider the private signal structure ρ^c , which randomizes with equal probabilities between two perfectly informative signal functions

$$\begin{aligned} \xi_1(\theta) &= \theta, \text{ and} \\ \xi_2(\theta) &= \begin{cases} \theta & \text{if } \theta < \theta_0, \\ 1 + \theta_0 - \theta & \text{if } \theta \geq \theta_0. \end{cases} \end{aligned} \quad (21)$$

Fig. 2 depicts signal functions $\xi_i(\theta)$, $i = 1, 2$. The signal sets, that is, the images of ξ_1 and ξ_2 are identical and equal to $\mathbf{S} = [0, 1]$. Thus, any buyer's deviation from truth-telling

³⁵The HARA function is defined as $v(x, \gamma) = \frac{1-\gamma}{\gamma} \left(\frac{\alpha}{1-\gamma}x + \beta \right)^\gamma$. Adding the constant term $-\frac{1-\gamma}{\gamma} \beta^\gamma$ is a normalization, which does not affect the shape of the function, but guarantees that condition (16) holds.

³⁶See Ingersoll (1987).

is undetectable by the seller. Also, the buyer perfectly infers $\theta = s$ upon observing a signal $s < \theta_0$, but is uncertain about θ upon observing $s \geq \theta_0$. In the latter case, the posterior probability $q(\theta|s, \rho^c)$ of state θ induced by signal s under a private signal structure ρ^c is the binary distribution, which places probabilities $\frac{1}{2}$ on the posterior states

$$\begin{aligned}\theta_1 &= \varphi_1(s) = s, \text{ and} \\ \theta_2 &= \varphi_2(s) = \begin{cases} s & \text{if } s < \theta_0, \\ 1 + \theta_0 - s & \text{if } s \geq \theta_0. \end{cases}\end{aligned}$$

where $\varphi_i = \xi_i^{-1}$, $i = 1, 2$.

Also, consider the direct mechanism $\mathcal{M}^c = (Q_{\xi_i}^c(s), t_{\xi_i}^c(s))$, $i = 1, 2$ with the message set $\mathbf{M} = \mathbf{S} = [0, 1]$, such that

$$\begin{aligned}Q_{\xi_i}^c(s) &= Q^c(s) = \begin{cases} 0 & \text{if } s < \theta_0, \\ 1 & \text{if } s \geq \theta_0, \end{cases} \\ t_{\xi_i}^c(s) &= \begin{cases} 0 & \text{if } s < \theta_0, \\ \varphi_i(s) & \text{if } s \geq \theta_0. \end{cases}\end{aligned}$$

Given the pair (ρ^c, \mathcal{M}^c) , the buyer's posterior payoff is

$$\begin{aligned}EV_B(m|s, \gamma, \rho^c) &= 0 \text{ if } m < \theta_0, s \in \mathbf{S}, \text{ and} \\ EV_B(m|s, \gamma, \rho^c) &= q_1(s, \rho^c) V(t_{\xi_1}^c(m), \varphi_1(s), \gamma) + q_2(s, \rho^c) V(t_{\xi_2}^c(m), \varphi_2(s), \gamma) \\ &= \frac{1}{2}v(\varphi_1(s) - \varphi_1(m), \gamma) + \frac{1}{2}v(\varphi_2(s) - \varphi_2(m), \gamma) \\ &= \begin{cases} \frac{1}{2}v(s - m, \gamma) + \frac{1}{2}v(m - s, \gamma) & \text{if } m \geq \theta_0, s \geq \theta_0, \\ \frac{1}{2}v(s - m, \gamma) + \frac{1}{2}v(s - (1 + \theta_0 - m), \gamma) & \text{if } m \geq \theta_0, s < \theta_0. \end{cases}\end{aligned}$$

It is straightforward to show that $EV_B(m|s, \gamma, \rho^c)$ is maximized at $m = s$ for all $(s, \gamma) \in \mathbf{S} \times \mathbf{T}$, and $EV_B(s|s, \gamma, \rho^c) = 0$ for all $s \in \mathbf{S}$. Hence, the interim incentive-compatibility constraints (19) and the interim individual-rationality constraints (20) hold. Furthermore, the buyer's ex-post payoff for $\theta_i = \varphi_i(s) \geq \theta_0$ is

$$\begin{aligned}V(\varphi_i(s), t_{\xi_i}^c(s), \gamma) &= v(\varphi_i(s) - t_{\xi_i}^c(s), \gamma) \\ &= v(\varphi_i(s) - \varphi_i(s), \gamma) = 0 \text{ if } s \geq \theta_0, \gamma \in \mathbf{T}, i = 1, 2.\end{aligned}$$

This implies that the seller extracts the full surplus in each state $\theta \geq \theta_0$ for any type $\gamma \in \mathbf{T}$ upon inferring θ from m and ξ_i .

Intuitively, the possibility for the seller to extract the full information and surplus from the buyer without violating her interim incentive-compatibility and individual-rationality constraints is driven by a combination of three factors. First, the object is sold to the buyer if and only if the mechanism infers that the state is above the cutoff θ_0 . That is, the object is allocated to the buyer if and only if her ex-post highest acceptable payment exceeds the seller's ex-post benefits from keeping the object. Second, the incentive-compatibility in these state is sustained by the opposite reactions of the buyer's payments under different signal functions. That is, any buyer's deviation in an attempt to reduce the payment under one

signal function is offset by the larger payment under the other signal function. A proper selection of φ_1 and φ_2 eliminates the buyer's marginal benefits from both local distortions (that is, when $s \geq \theta_0$ and $m \geq \theta_0$) and global ones (that is, when $s < \theta_0$ and $m \geq \theta_0$) and thus sustains buyer's truth-telling.³⁷ As a result, the mechanism perfectly infers the posterior state θ_i from m and the realized signal function ξ_i . Third, the above effect does not depend on the absolute values of buyers' payments. Thus, the seller can charge the buyer with the maximum payment, which precludes her from selecting the outside option. Because the value of this payment does not depend on the buyer's type γ , the mechanism extracts the full surplus from the buyer in all target states.

Two comments are worth mentioning. First, similarly to Example A, the full surplus extraction is feasible regardless of the precision of the buyer's information. In particular, if $s = \frac{1+\theta_0}{2}$, then the buyer is perfectly informed about the posterior state $\hat{\theta} = \hat{s} = \frac{1+\theta_0}{2}$. However, she is still unable to receive a positive payoff by using this information.

Second, the analysis above does not require the strict concavity of $V(t, \theta, \gamma)$ in t . Hence, it is equally applicable to the buyer's payoff function, which is linear in θ and t for all γ .³⁸

$$V(t, \theta, \gamma) = v(\theta - t, \gamma) = \alpha(\gamma)(\theta - t).$$

Because the buyer is risk-neutral in this case, her interim payoff EV_B is unaffected by lotteries over payments under different signal functions, which are induced by her message. In other words, the buyer is indifferent between all messages upon receiving a signal $s \geq \theta_0$:

$$EV_B(m|s, \gamma, \rho^c) = \frac{1}{2}v(s - m, \gamma) + \frac{1}{2}v(m - s, \gamma) = 0 \text{ for } m \in \mathbf{S}, s \geq \theta_0.$$

4.3 Full surplus extraction

In this subsection we establish the possibility of the full information and surplus extraction for states above an arbitrary cutoff $\theta_0 \in \Theta$ in the general case. Consider the private signal structure ρ^* , which randomizes with equal probabilities between signal functions $\xi_1 = \varphi_1^{-1} : \Theta \rightarrow \mathbf{S}$ and $\xi_2 = \varphi_2^{-1} : \Theta \rightarrow \mathbf{S}$, where $\mathbf{S} = [\underline{s}, \bar{s}]$, $\varphi_1' > 0$, and

$$\varphi_2(s) = \begin{cases} \varphi_1(s) & \text{if } s < s_0, \\ \Psi^{-1}(\Psi(\theta_0) + \Psi(\bar{\theta}) - \Psi(\varphi_1(s))) & \text{if } s \geq s_0, \end{cases} \quad (22)$$

where $s_0 = \xi_1(\theta_0)$, or equivalently, $\varphi_1(s_0) = \theta_0$. Thus, upon receiving a signal $s < s_0$ the buyer perfectly infers the state $\theta = \varphi_1(s)$. For $s \geq s_0$, the buyer's posterior belief is a binary distribution over $\{\varphi_1(s), \varphi_2(s)\}$, where $\varphi_2(s)$ satisfies the differential equation (14) with the boundary condition $\varphi_2(s_0) = \bar{\theta}$.

³⁷Verifying that \mathcal{M} is interim incentive-compatible for other combinations of s and m is trivial.

³⁸In general, $v(\theta - t, \gamma) = \alpha(\gamma)(\theta - t) + \beta(\gamma)$. However, condition (16) implies $\beta(\gamma) = 0$.

Next, consider a mechanism $\mathcal{M}^{\rho^*} = \left(Q_{\xi_i}^{\rho^*}(s), t_{\xi_i}^{\rho^*}(s) \right)$, such that

$$Q_{\xi_i}^{\rho^*}(s) = Q^{\rho^*}(s) = \begin{cases} 0 & \text{if } s < s_0, \\ 1 & \text{if } s \geq s_0, \end{cases} \quad (23)$$

$$t_{\xi_i}^{\rho^*}(s) = \begin{cases} 0 & \text{if } s < s_0, \\ \varphi_i(s) & \text{if } s \geq s_0. \end{cases} \quad (24)$$

The theorem below establishes that the pair $(\rho^*, \mathcal{M}^{\rho^*})$ extracts the full information and surplus from the buyer for states $\theta \geq \theta_0$ under Condition 1.

Theorem 2 *Suppose V satisfies (16)–(18) and (f, V) satisfy Condition 1 for Θ_0 . Consider the private signal structure ρ^* that randomizes between $\xi_1 = \varphi_1^{-1}$ and $\xi_2 = \varphi_2^{-1}$ with equal probabilities, where $\varphi_1 : \mathbf{S} \rightarrow \Theta$ is differentiable, $\varphi_1' > 0$, and φ_2 is given by (22). Then ρ^* and the mechanism $\mathcal{M}^{\rho^*} = \left(Q_{\xi_i}^{\rho^*}(s), t_{\xi_i}^{\rho^*}(s) \right)$ extract the full surplus for $(\theta, \gamma) \in \Theta_0 \times \mathbf{T}$.*

The proof of theorem consists of two parts. The first part demonstrates that the mechanism \mathcal{M}^{ρ^*} is interim individually-rational under the private signal structure ρ^* . Specifically, for signals $s < s_0$, the buyer receives the outside option with value 0. For $s \geq s_0$, the buyer pays θ_i in each posterior state $\theta_i = \varphi_i(s)$, $i = 1, 2$, which is equal to her willingness to pay. That is, the mechanism extracts the buyer's full surplus upon learning the state θ from $m = s$ and ξ_i .

The main part of the proof is to establish the interim incentive-compatibility of the mechanism \mathcal{M}^{ρ^*} , which is done in a few steps depending on the values of a signal s and a message m . For $s \geq s_0$ and $m \geq s_0$, the interim incentive-compatibility is an implication of Theorem 1. First, the incentive-compatibility for these values of s and m is equivalent to the optimality of the truthful strategy in the implementation model with the state space Θ_0 , the prior density $f_0(\theta) = f(\theta|\theta \in \Theta_0)$, the signal set \mathbf{S}_0 , and the private signal structure ρ_0 that randomizes between $\xi_1(\theta)$ and $\xi_2(\theta)$ with the domains restricted to Θ_0 . Second, the inverses φ_1 and φ_2 satisfy the first-order condition (14) with the boundary condition $\varphi_1(s_0) = \theta_0$ in the equivalent implementation model. Third, because V and f satisfy (17) and Condition 1 for Θ_0 , then applying Theorem 1 to the equivalent implementation model means that the agent's truthful strategy is optimal. This in turn results in the interim incentive-compatibility of \mathcal{M}^{ρ^*} for $(s, m) \in \mathbf{S}_0^2$. Next, for $m < s_0$ the incentive-compatibility holds as the buyer receives her outside option of value 0 for all $s \in \mathbf{S}$, which is identical to her payoff from truthful reporting. This is because truthful reporting provides the buyer with the outside option for $s < s_0$. For $s \geq s_0$, truthful reporting results in the full surplus extraction, so the buyer receives 0 as well. The final step is to show that the buyer with a signal $s < \theta_0$ cannot benefit from deviating to $m \geq s_0$. This step is based on the monotonicity of the buyer's payoff $V(t, \theta, \gamma)$ in θ and the fact that the buyer with signal $s \geq s_0$ does not receive a positive surplus. This completes the proof of the theorem.

Importantly, conditions (16) and (18) are essential for the full surplus extraction. Without additional assumptions about the impact of buyer's private information on her preferences, the seller cannot extract the full surplus by using the private information design even in the simplest model with quasi-linear players' preferences and discrete states.³⁹ Importantly,

³⁹See Remark 8 in Krämer (2020).

these conditions are local. This is because for a given state θ , they must hold only at the ‘full surplus extraction’ point $t = \theta$, i.e., for the payment equal to the buyer’s willingness to pay. Equivalently, they must hold only along the diagonal (θ, θ) in the (t, θ) space.

5 Conclusion and discussion

This paper adds to the literature on the agency problem by showing how the principal can use private information design in a simple way to implement her ideal action for a target subset of states or the entire state space. The result holds even if the agent’s preferences are non-quasilinear, non-concave, depend on the privately known component, and are independent of the state.

We conclude the paper by suggesting possible avenues for future research. First, the proposed construction of private signal structures can be potentially used in other economic environments. These may include models in which players’ ideal actions are non-monotone to the unknown information or the buyer’s payoff is non-monotone in the principal’s action. Intuitively, truthtelling of the agent is driven by opposite monotonicities of the payoffs in her message for different posterior states. In general, each of these payoffs is a composition of three functions: i) the payoff as a function of the principal’s action; ii) the principal’s ideal action as a function of the posterior state; and iii) the induced posterior state as a function of the agent’s message, which is the inverse of the signal function.⁴⁰ In our paper, we assume the strict monotonicity of the first two functions. If one or both of these functions are non-monotone, then the monotonicity of the composite function can be potentially restored by selecting a non-monotone (but bijective and, hence, perfectly informative) inverse.

Another avenue for future research is to extend the setup to multidimensional state and action spaces. If the agent’s payoff function is additively separable, then our construction can be easily applied coordinatewise.⁴¹ However, the question of whether our construction can be extended to multidimensional spaces in the case of payoff functions of the general form remains open.

Appendix

Proof of Theorem 1 Consider functions $\varphi_i : \mathbf{S} \rightarrow \Theta, i = 1, 2$, such that φ_1 is differentiable, $\varphi_1' > 0$, and φ_2 is given by (15). By construction, the pair $\{\varphi_1, \varphi_2\}$ satisfies the first-order condition (12). Because $\varphi_2 : \mathbf{S} \rightarrow \Theta$ is such that $\varphi_2' < 0$, then it is bijective. Hence, the functions $\xi_i = \varphi_i^{-1} : \Theta \rightarrow \mathbf{S}, i = 1, 2$ exist and are perfectly informative signal functions. Consider the private signal structure ρ^* , which randomizes between ξ_1 and ξ_2 with equal probabilities.

⁴⁰The agent’s posterior payoff is $EV(m|s, \gamma, \rho) = \sum_{i=1}^2 q_i(s, \rho) V(y(\varphi_i(m)), \varphi_i(s), \gamma)$. Hence, the payoff conditional on the posterior state $\theta_i = \varphi_i(s)$ is given by $V(y(\varphi_i(m)), \varphi_i(s), \gamma)$.

⁴¹Consider, for instance, $V(\vec{a}, \vec{\theta}) = -\sum_{i=1}^2 (a_i - \theta_i - b_i)^2$, where $\vec{a} = (a_1, a_2), \vec{\theta} = (\theta_1, \theta_2)$, and $\vec{\theta}$ is uniformly distributed on $[0, 1]^2$. Then the private signal structure, which randomizes between signal functions $\xi_1(\vec{\theta}) = \vec{\theta}$ and $\xi_2(\vec{\theta}) = (1, 1) - \vec{\theta}$ with equal probabilities, sustains agent’s truthtelling.

Next, the truthful strategy is optimal for the agent if $EV(m|s, \gamma, \rho^*)$ given by (9) is pseudo-concave in m for all $(s, \gamma) \in \Theta \times \mathbf{T}$.⁴² To establish the pseudo-concavity of $EV(m|s, \gamma, \rho^*)$ in m , it is sufficient to show that the function

$$\phi(m|s, \gamma, \rho^*) = \frac{\partial}{\partial m} EV(m|s, \gamma, \rho^*)$$

is pseudo-monotone in m on \mathbf{S} for all $(s, \gamma) \in \mathbf{S} \times \mathbf{T}$ (Proposition 2.5, Hadjisavvas et al. 2005). A function $\phi(m|s, \gamma, \rho^*)$ is pseudo-monotone in m on an interval $\mathbf{S} \subset \mathbb{R}$ if $\phi(m_1|s, \gamma, \rho^*)(m_2 - m_1) \leq 0$ implies $\phi(m_2|s, \gamma, \rho^*)(m_2 - m_1) \leq 0$ for all $(m_1, m_2) \in \mathbf{S}^2$.

To guarantee the pseudo-monotonicity of $\phi(m|s, \gamma, \rho^*)$, we use the aggregation result by Quah and Strulovici (Proposition 1, 2012). It says that a linear combination $\alpha_1 \mathcal{V}_1(m) + \alpha_2 \mathcal{V}_2(m)$ of two pseudo-monotone functions $\mathcal{V}_1(m)$ and $\mathcal{V}_2(m)$ is pseudo-monotone for all $\alpha_i \geq 0, i = 1, 2$ if and only if: (i) $-\frac{\mathcal{V}_1(m)}{\mathcal{V}_2(m)}$ is decreasing in m for all m such that $\mathcal{V}_1(m) > 0$ and $\mathcal{V}_2(m) < 0$; and (ii) $-\frac{\mathcal{V}_2(m)}{\mathcal{V}_1(m)}$ is decreasing in m for all m such that $\mathcal{V}_1(m) < 0$ and $\mathcal{V}_2(m) > 0$.⁴³

Fix $(s, \gamma) \in \mathbf{S} \times \mathbf{T}$. It follows from (12) that

$$\phi(m|s, \gamma, \rho^*) = \sum_{i=1}^2 q_i(s, \rho^*) V'_a(\varphi_i(m), \varphi_i(s), \gamma) \varphi'_i(m) = \sum_{i=1}^2 q_i(s, \rho^*) V_i(m, s, \gamma),$$

where

$$\mathcal{V}_i(m, s, \gamma) = V'_a(\varphi_i(m), \varphi_i(s), \gamma) \varphi'_i(m), i = 1, 2.$$

Then, $V'_a < 0, \varphi'_1 > 0$, and $\varphi'_2 < 0$ imply

$$\mathcal{V}_1(m, s, \gamma) < 0 \text{ and } \mathcal{V}_2(m, s, \gamma) > 0 \text{ for all } (m, s, \gamma) \in \mathbf{S}^2 \times \mathbf{T}.$$

Therefore, $\mathcal{V}_i(m, s, \gamma) < 0, i = 1, 2$ is strictly pseudo-monotone in m for all $(m, s, \gamma) \in \mathbf{S}^2 \times \mathbf{T}$. Thus, $\phi(m|s, \gamma, \rho^*)$ is pseudo-monotone in m for $q_i(s, \rho^*) \geq 0, i = 1, 2$ if and only if

$$-\frac{\mathcal{V}_2(m, s, \gamma)}{\mathcal{V}_1(m, s, \gamma)} = -\frac{V'_a(\varphi_2(m), \varphi_2(s), \gamma) \varphi'_2(m)}{V'_a(\varphi_1(m), \varphi_1(s), \gamma) \varphi'_1(m)}$$

is decreasing in m . By using (14), we get

$$-\frac{\varphi'_2(s)}{\varphi'_1(s)} = \frac{h(\varphi_1(s))}{h(\varphi_2(s))},$$

⁴²A differentiable function $\mathcal{V}(x)$ is *pseudo-concave* on a convex set $\mathbf{X} \subset \mathbb{R}$ if for every $(x, y) \in \mathbf{X}^2$, $\mathcal{V}(x) < \mathcal{V}(y)$ implies $\mathcal{V}'(x)(y - x) > 0$. If $\mathcal{V}'(x_0) = 0$ for $x_0 \in \mathbf{X}$, then x_0 is a maximizer of \mathcal{V} (Proposition 2.4 in Hadjisavvas et al., 2005).

⁴³Quah and Strulovici (2012) use the term *single crossing ϕ* , which is equivalent to a pseudo-monotone $-\phi$. Formally, a single crossing function can intersect the x -axis at a single interval from below, whereas a pseudo-monotone function can intersect the x -axis at a single interval from above.

and

$$\begin{aligned} -\frac{\mathcal{V}_2(m, s, \gamma)}{\mathcal{V}_1(m, s, \gamma)} &= \frac{V'_a(\varphi_2(m), \varphi_2(s), \gamma) \sqrt{-f(\varphi_1(m))\zeta(\varphi_1(m))}}{V'_a(\varphi_1(m), \varphi_1(s), \gamma) \sqrt{-f(\varphi_2(m))\zeta(\varphi_2(m))}} \\ &= \frac{\nu(\varphi_2(m), \varphi_2(s), \gamma)}{\nu(\varphi_1(m), \varphi_1(s), \gamma)}, \end{aligned}$$

where

$$\nu(a, \theta, \gamma) = \frac{V'_a(a, \theta, \gamma)}{\sqrt{-f(a)\zeta(a)}} = \frac{V'_a(a, \theta, \gamma)}{h(a)}.$$

Because $V'_a < 0$ and $h > 0$, it follows that $\nu(a, \theta, \gamma) < 0$ for all $(a, \theta, \gamma) \in \Theta^2 \times \mathbf{T}$. Also, $\varphi_i \in \Theta, i = 1, 2$. Then, $\varphi'_1 > 0$ and Condition 1 imply that $\nu(\varphi_1(m), \varphi_1(s), \gamma)$ is decreasing in m for all $(s, \gamma) \in \mathbf{S} \times \mathbf{T}$. Similarly, $\varphi'_2 < 0$ and Condition 1 imply that $\nu(\varphi_2(m), \varphi_2(s), \gamma)$ is increasing in m for all $(s, \gamma) \in \mathbf{S} \times \mathbf{T}$. By combining these arguments, it follows that $-\frac{\mathcal{V}_1(m, s, \gamma)}{\mathcal{V}_2(m, s, \gamma)}$ is decreasing in m for all $(s, \gamma) \in \mathbf{S} \times \mathbf{T}$. ■

Proof of Theorem 2 Consider functions $\varphi_i : \mathbf{S} \rightarrow \Theta, i = 1, 2$, such that φ_1 is differentiable, $\varphi'_1 > 0$, and φ_2 is given by (22). Because $\varphi_2 : \mathbf{S} \rightarrow \Theta$ is piecewise continuous and strictly monotone for $s < \theta_0$ and $s \geq \theta_0$, and $\varphi_2(s) < \theta_0 \leq \varphi_2(z)$ for all $s < \theta_0 \leq z$, then φ_2 is bijective. Hence, the functions $\xi_i = \varphi_i^{-1} : \Theta \rightarrow \mathbf{S}, i = 1, 2$ exist and are perfectly informative signal functions. Next, consider the private signal structure ρ^* , which randomizes with equal probabilities between ξ_1 and ξ_2 , and the mechanism \mathcal{M}^{ρ^*} with the allocation and payment rules given by (23) and (24), respectively.

First, the interim individual-rationality constraints (20) hold and are binding for all $s \in \mathbf{S}$. If $s < s_0$, then $Q_{\xi_i}^{\rho^*}(s) = 0$ and $t_{\xi_i}^{\rho^*}(s) = 0, i = 1, 2$. That is, the buyer receives the outside option, and her posterior payoff $EV_B(s|s, \gamma, \rho) = V(0, 0, \gamma) = 0$ for all $s < s_0$ and $\gamma \in \mathbf{T}$. If $s \in \mathbf{S}_0 = [s_0, \bar{s}]$, then $Q_{\xi_i}^{\rho^*}(s) = 1$ and $t_{\xi_i}^{\rho^*}(s) = \varphi_i(s), i = 1, 2$. This results in

$$EV_B(s|s, \gamma, \rho^*) = \sum_{i=1}^2 q_i(s, \rho^*) V(\varphi_i(s), \varphi_i(s), \gamma) = 0 \text{ for all } (s, \gamma) \in \mathbf{S}_0 \times \mathbf{T}. \quad (25)$$

where the second equality holds because (16) implies $V(\varphi_i(s), \varphi_i(s), \gamma) = 0$ for $(s, \gamma) \in \Theta \times \mathbf{T}, i = 1, 2$.

Second, we prove the interim incentive-compatibility of the pair $(\rho^*, \mathcal{M}^{\rho^*})$ by considering three cases depending on the values of $(m, s) \in \mathbf{S}^2$.

(i) $(s, m) \in \mathbf{S}_0^2$. Because $Q_{\xi_i}^{\rho^*}(m) = 1$ and $t_{\xi_i}^{\rho^*}(m) = \varphi_i(m), i = 1, 2$ for $m \in \mathbf{S}_0$, we have

$$EV_B(m|s, \gamma, \rho^*) = \sum_{i=1}^2 q_i(s, \rho^*) V(\varphi_i(m), \varphi_i(s), \gamma) \text{ for } (m, s, \gamma) \in \mathbf{S}_0^2 \times \mathbf{T}. \quad (26)$$

Next, note that states $\theta \in \Theta_0$ generate signals $s_i = \varphi_i^{-1}(\theta) \in \mathbf{S}_0, i = 1, 2$ under ρ^* . Hence, the agent's posterior belief induced by a signal $s \in \mathbf{S}_0$ is the binary distribution over states

$\theta_i = \varphi_i(s) \in \Theta_0, i = 1, 2$. By using (7), the posterior probability of θ_i is

$$\begin{aligned} q_i(s, \rho^*) &= \frac{f(\varphi_i(s))\varphi'_i(s)}{f(\varphi_1(s))\varphi'_1(s) - f(\varphi_2(s))\varphi'_2(s)} \\ &= \frac{f_0(\varphi_i(s))\varphi'_i(s)}{f_0(\varphi_1(s))\varphi'_1(s) - f_0(\varphi_2(s))\varphi'_2(s)} = q_i(s, \rho^0), i = 1, 2. \end{aligned}$$

Here, $f_0(\theta) = f(\theta|\theta \in \Theta_0) = \frac{f(\theta)}{1-F(\theta_0)}$ is the prior density of θ conditional on $\theta \in \Theta_0$, $F(\theta)$ is the cdf of θ , and ρ^0 is the private signal function that randomizes with equal probabilities between ξ_1^0 and ξ_2^0 , where $\xi_i^0 = \xi_i : \Theta_0 \rightarrow \mathbf{S}_0$ is a signal function ξ_i with the domain restricted by Θ_0 and, thus, the image \mathbf{S}_0 .

Now, consider the implementation model with the signal set \mathbf{S}_0 , the prior density $f_0(\theta)$, and the private signal structure ρ^0 . Hereafter, we call this model the *equivalent implementation model*. By combining the arguments above and comparing (26) with (9), it follows that the interim incentive-compatibility condition (19) for $(s, \gamma) \in \mathbf{S}_0 \times \mathbf{T}$ on the restricted message space \mathbf{S}_0 is identical to the optimality condition (11) for the agent's truthful strategy in the equivalent implementation model. Next, $\varphi_2(\theta)$ given by (22) satisfies the first-order condition (14) with the boundary condition $\varphi_1(s_0) = \theta_0$ in this model. Also, conditions (17) and 1 hold for Θ_0 . Then, by applying Theorem 1 to the equivalent implementation model, it follows that the agent's truthful strategy is optimal. This means that the incentive-compatibility constraints in the bilateral-trade model hold for $(s, m) \in \mathbf{S}_0^2$.

(ii) $s \in \mathbf{S}, m < s_0$. Then $Q_{\xi_i}^{\rho^*}(m) = 0, t_{\xi_i}^{\rho^*}(m) = 0, i = 1, 2$, and

$$EV_B(m|s, \gamma, \rho^*) = V(0, 0, \gamma) = 0 = EV_B(s|s, \gamma, \rho^*),$$

where $EV_B(s|s, \gamma, \rho^*) = 0$ for $(s, \gamma) \in \mathbf{S} \times \mathbf{T}$ follows from the interim individual-rationality of \mathcal{M}^* .

(iii) $s < s_0, m \geq s_0$. Since $s < s_0$, then $\varphi_i(s) < \theta_0 = \varphi_1(s_0), i = 1, 2$. Also, $m \geq s_0$ implies $Q_{\xi_i}^{\rho^*}(m) = 1$ and $t_{\xi_i}^{\rho^*}(m) = \varphi_i(s), i = 1, 2$. Then, we have

$$\begin{aligned} EV_B(m|s, \gamma, \rho^*) &= \sum_{i=1}^2 q_i(s, \rho^*) V(\varphi_i(m), \varphi_i(s), \gamma) < \sum_{i=1}^2 q_i(s, \rho^*) V(\varphi_i(m), \varphi_1(s_0), \gamma) \\ &\leq \sum_{i=1}^2 q_i(s, \rho^*) V(\varphi_i(m), \varphi_i(s_0), \gamma) \leq \sum_{i=1}^2 q_i(s, \rho^*) V(\varphi_i(s_0), \varphi_i(s_0), \gamma) \\ &= EV_B(s_0|s_0, \gamma, \rho^*) = 0, \end{aligned}$$

where the first inequality holds since $\varphi_i(s) < \varphi_1(s_0), i = 1, 2$ and $V'_\theta(t, \theta, \gamma) > 0$ imply $V(\varphi_i(m), \varphi_i(s), \gamma) < V(\varphi_i(m), \varphi_1(s_0), \gamma), i = 1, 2$, the second inequality holds due to $\varphi_2(s_0) = \bar{\theta} \geq \theta_0 = \varphi_1(s_0)$ and $V'_\theta(t, \theta, \gamma) > 0$, and the last one holds since $m = s_0$ maximizes $EV_B(m|s_0, \gamma, \rho)$ over $m \geq s_0$. ■

References

- Bergemann, D. and M. Pesendorfer, 2007. Information structures in optimal auctions. *Journal of Economic Theory* 137, 580–609
- Bergemann, D., S. Morris, and T. Heumann, 2022. Screening with persuasion. Working paper
- Blume, A., Board, O., and K. Kawamura, 2007. Noisy talk. *Theoretical Economics* 2, 395–440
- Blume, A., Lai, E., and W. Lim, 2019. Eliciting private information with noise: the case of randomized response. *Games and Economic Behavior* 113, 356–380
- Chakraborty, A. and R. Harbaugh, 2010. Persuasion by cheap talk. *American Economic Review* 100, 2361–2382
- Crawford, V., and J. Sobel, 1982. Strategic information transmission. *Econometrica* 50, 1431–1451
- Crémer, J., and R.P. McLean, 1988. Full extraction of the surplus in Bayesian and dominant strategy auctions. *Econometrica* 56, 1247–1257
- Dworzak, P., Kominers, S.D., and M. Akbarpour, 2021. Redistribution through markets. *Econometrica* 89, 1665–1698
- Ederer, F., Holden, R., and M. Meyer, 2018. Gaming and strategic opacity in incentive provision. *RAND Journal of Economics* 49, 819–854
- Esö, P., and B. Szentes, 2007. Optimal information disclosure in auctions. *Review of Economic Studies* 74, 705–731
- Fu, H., Haghpanah, N., Hartline, J., and R. Kleinberg, 2021. The full surplus extraction from samples. *Journal of Economic Theory* 193, 105230
- Goltsman, M., Hörner, J., Pavlov, G., and F. Squintani, 2009. Mediation, arbitration and negotiation. *Journal of Economic Theory* 144, 1397–1420
- Hadjisavvas, N., Komlòsi, S., and S. Schaible, 2005. *Handbook of generalized convexity and generalized monotonicity*. Springer, Boston
- Heifetz, A., and Z. Neeman, 2006. On the generic (im)possibility of full surplus extraction in mechanism design. *Econometrica* 74, 213–233
- Hwang, I., Kim, T., and R. Boleslavsky, 2019. Competitive advertising and pricing. Working paper
- Ingersoll, J., 1987. *Theory of financial decision making*. Totowa, NJ: Rowman and Littlefield Publishers
- Ivanov, M. and A. Sam, 2022. Cheap talk with private signal structures, *Games and Economic Behavior* 132, 288–304
- Ivanov, M., 2021. Optimal monotone signals in Bayesian persuasion mechanisms. *Economic Theory* 72, 955–1000
- Ivanov, M., 2016. Dynamic learning and strategic communication. *International Journal of Game Theory* 45, 627–653
- Ivanov, M., 2015. Dynamic information revelation in cheap talk. *The B.E. Journal of Theoretical Economics* 15, 251–275
- Ivanov, M., 2013. Information revelation in competitive markets. *Economic Theory* 52, 337–365

- Ivanov, M., 2010. Communication via a strategic mediator. *Journal of Economic Theory* 145, 869–884
- Johnson, J., and D. Myatt, 2006. On the simple economics of advertising, marketing, and product design. *American Economic Review* 93, 756–784
- Kamenica, E. and M. Gentzkow, 2011. Bayesian persuasion. *American Economic Review* 101, 2590–2615
- Krähmer, D., 2020. Information disclosure and the full surplus extraction in mechanism design. *Journal of Economic Theory* 187, 105020
- Krähmer, D., 2021. Information design and strategic communication. *American Economic Review: Insights* 3, 51–66
- Krishna, V., 2009. *Auction Theory*, 2nd edition. Elsevier Academic Press, Boston
- Larionov, D., Pham, H., and T. Yamashita, 2021. First best implementation with costly information acquisition, Working paper
- Lewis, T. and D. Sappington, 1994. Supplying information to facilitate price discrimination. *International Economic Review* 35, 309–327
- Li, H., and X. Shi, 2019. Discriminatory information disclosure. *American Economic Review* 107, 3363–3385
- Lipnowski, E. and D. Ravid, 2020. Cheap talk with transparent motives. *Econometrica* 88, 1631–1660
- McAfee, P. and P. Reny, 1992. Correlated information and mechanism design. *Econometrica* 60, 395–421
- McKinsey & Company. How companies spend their money: a McKinsey global survey, *McKinsey Quarterly*, June 2007.
- Myerson, R., 1981. Optimal auction design. *Mathematics of Operations Research* 6, 58–73
- Pastrian, N., 2021. Surplus extraction with behavioral types, Working paper
- Quah, J. and B. Strulovici, 2012. Aggregating the single crossing property. *Econometrica* 80, 2333–2348
- Saak, A., 2006. The optimal private information in single unit monopoly. *Economics Letters* 91, 267–272
- Shannon, C., 1949. Communication theory of secrecy systems. *Bell System Technical Journal* 28, 656–715
- Watson, J., 1996. Information transmission when the informed party is confused. *Games and Economic Behavior* 12, 143–161
- Zhu, S., 2021. Private disclosure with multiple agents, Working paper.