

Acquisition of/Stochastic Evidence¹

Elchanan Ben-Porath ² Eddie Dekel ³ Barton L. Lipman⁴

Preliminary Draft
October 2019

¹We thank the National Science Foundation, grant SES-1919319 (Dekel and Lipman), and the US-Israel Binational Science Foundation for support for this research.

²Department of Economics and Center for Rationality, Hebrew University. Email: benporat@math.huji.ac.il.

³Economics Department, Northwestern University, and School of Economics, Tel Aviv University. Email: dekel@northwestern.edu.

⁴Department of Economics, Boston University. Email: blipman@bu.edu.

Abstract

We explore two highly interrelated models of “hard information.” In the *evidence–acquisition model*, an agent with private information searches for evidence to show to the principal about her type. In the *signal–choice model*, a privately informed agent chooses an action which generates a random signal whose realization may be correlated with her type. We show that the signal–choice model is a special case and, under certain conditions, a reduced form of the evidence–acquisition model. We develop tools for characterizing optimal mechanisms for these models by giving conditions under which some aspects of the principal’s optimal choices can be identified only from the information structure, without regard to the utility functions or the principal’s priors. We also give a novel result on conditions under which there is no value to commitment for the principal.

1 Introduction

The earliest work in economics on information transmission considered settings where any agent, regardless of her information, could send any “signal” or “message,” though potentially at costs which depend on her private information. In Spence’s (1973) classic signaling model, the cost of education depends on the agent’s type, but she can get any amount of education she wants. In Crawford and Sobel’s (1982) classic cheap-talk model, there are no costs, so nothing about the messages depend on the agent’s type. In such models, effective communication relies on how the agent’s incentives to take different actions or to induce different beliefs varies with her type.¹

More recent work treats information as more strongly tied to the agent’s type. For example, in the literature on hard evidence, the set of messages the agent has available can depend on her type, so that presentation of certain messages proves certain facts about her type. In the literature on career concerns, the agent’s actions affect observable outcomes differently for different types. In the literature on Bayesian persuasion, the outcome of a given experiment carried out by the agent depends on her type.

We connect these and other approaches to information transmission through a model where the agent chooses among actions that generate random signals depending on her type. The agent can then choose which realizations to present to a principal who chooses an action affecting the utility of both the principal and the agent. We refer to this as the *evidence-acquisition model*. We also study a model that is a special case and, under some conditions, a reduced form of the evidence-acquisition model, called the *signal-choice model* in which the agent has no option about what evidence to present. In addition to providing a unification of this literature, these models raise new issues in game theory and mechanism design, giving a new approach to analyzing information transmission.

In Section 2, we present the “technology” of the two models, relate them to the literature, and introduce a running example. In Section 3, we briefly discuss game-theoretic versions of the models and show that in a natural game, the signal-choice model is a reduced form of the evidence-acquisition model.

In Section 4, we turn to mechanism design. First, we provide a useful analog of the Revelation Principle for the evidence-acquisition model. The general class of mechanisms for these problems is quite complex, involving numerous steps of communication between the agent and the principal. We analyze conditions under which we can identify some of the principal’s communication in the optimal mechanism, specifically, the principal’s recommendation regarding what evidence he would like to see. When we can identify

¹This description omits the literature on cheap talk with type-independent preferences for the sender, largely initiated by Chakraborty and Harbaugh (2010). See Lipnowski and Ravid (2019) for a recent example.

this recommendation, we can reduce the model to the signal–choice model. We also give conditions under which we can identify the principal’s recommendation regarding what signal to choose, leading to a further simplification.

In Section 5, we show that under certain conditions, the optimal mechanism does not require commitment by the principal. That is, the best mechanism for the principal yields the same outcome as the best equilibrium of the game where the principal is not committed. This result can be thought of as a generalization of Ben-Porath, Dekel, and Lipman (2019) or earlier results such as Glazer and Rubinstein (2004, 2006), Sher (2011), or Hart, Kremer, and Perry (2017). The result here is more general in that it allows stochastic evidence. This extension requires stronger assumptions on the preferences of the principal than those in our previous work. However, our assumptions on the preferences are weaker than those assumed in Glazer and Rubinstein (2004, 2006), and not comparable to Sher (2011) or Hart, Kremer, and Perry (2017). Our result applies in a wide class of interesting problems. In particular, for binary decisions such as, whether to hire a or not (at a fixed wage), whether or not to adopt a project, whether or not to provide some benefit or resource, and so on, the result is completely general in the sense that it applies for any utility functions for the agent and the principal.

Section 6 offers concluding remarks. Proofs not contained in the text are in the Appendix.

2 Models

In this section, we discuss the “primitives” of the model, reserving discussion of the specifics of the game or mechanism for later sections.

Running Example, Part 1. Throughout the paper, we will use the following example to illustrate ideas and results. We have an employer, also referred to as *the principal*, and an employee, also called *the agent*. The agent’s private information is her productivity for the principal. We consider two variations. First, we consider what we will call the *wage–setting version* of this problem. Here, as in Spence (1973), the principal sets a wage for the agent and his payoff is maximized by setting the wage equal to the agent’s true productivity. By contrast, the agent’s payoff is strictly increasing in the wage. Second, we consider the *hiring version* of the problem. Here there is a fixed wage, outside the control of the principal, and he can only decide whether or not to hire the agent. The principal prefers hiring to not hiring iff the agent’s productivity is sufficiently high, while the agent strictly prefers being hired, regardless of her true type. Hence in both versions the agent wants the principal to think she has a high type and the principal wants to know the true type, but in the second, the decision is coarser. We consider various

forms of evidence acquisition by the agent to try to persuade the principal she has high productivity.

As in the running example, the players in the model are an agent and a principal. The agent has a finite set of types T where the realization $t \in T$ is the agent’s private information. The principal’s prior over T is denoted τ and is assumed to have full support. The principal has a finite set of actions X . An element of X specifies all aspects of the principal’s action, including allocation of goods, monetary transfers, provision of public goods, or other activities. After information exchange between the agent and the principal, the principal chooses some $x \in X$. There is a set \mathcal{L} of all possible evidence messages. Information exchange includes the transmission of an evidence message and may also include cheap talk between the principal and the agent.

We consider two ways of modeling information transmission, one of which is a special case and, under certain conditions, a reduced form of the other. The more general model is the *evidence-acquisition model*, a model where the agent searches to find evidence. The agent has a variety of options available to try to obtain evidence. This search process could be sequential or one-shot. Rather than model this process, we focus on outcomes of the search process by treating the agent as choosing a probability distribution over the evidence set she ultimately obtains. We denote a typical set of evidence as $M \subseteq \mathcal{L}$ and a typical probability distribution over evidence sets as p_M . Let \mathcal{M}_t denote the set of such distributions that type t can generate. Thus $p_M \in \mathcal{M}_t \subseteq \Delta(2^{\mathcal{L}} \setminus \{\emptyset\})$.² Let \mathcal{M} be the set of possible message sets M that can be produced. That is, \mathcal{M} is the collection of M such that there exists t and $p_M \in \mathcal{M}_t$ with $M \in \text{supp}(p_M)$. The assumption that $\emptyset \notin \mathcal{M}$ can be thought of as assuming the agent can always say *something*, even if it is not informative — e.g., “I have no evidence to present.” If M is the realized set of messages, then the agent can present any one $m \in M$ to the principal.³ We write the utility function of the agent and principal as $u : T \times \mathcal{M} \times X \rightarrow \mathbf{R}$ and $v : T \times \mathcal{M} \times X \rightarrow \mathbf{R}$ respectively, where $\mathcal{M} = \cup_t \mathcal{M}_t$.

While we assume that the principal observes only the m sent by the agent and not the chosen evidence acquisition strategy p_M , the model allows observability of p_M as well. To see this, suppose that for every pair of distinct distributions, $p_M \neq p'_M$, if $m \in M \in \text{supp}(p_M)$, then there is no $M' \in \text{supp}(p'_M)$ with $m \in M'$. Then upon observing the message m , the principal knows the evidence acquisition strategy. Similarly, we can assume that certain distribution choices are observable and others are not or that only some of the messages reveal p_M in this sense, so that whether the distribution is observed is itself random and/or in the control of the agent.

²For any finite set B , $\Delta(B)$ is the set of probability distributions over B .

³As in the usual deterministic evidence model, the assumption that the agent can present only one message is without loss of generality. For example, if the agent could present two messages, we would simply replace \mathcal{L} with the set of pairs of messages.

Running Example, Part 2. For the purposes of the example, we assume a very stylized evidence–acquisition technology. To understand the idea, think of the agent of type t as able to choose a variety of ways to potentially demonstrate her ability. Each of these options gives a different probability distribution over an “outcome” she generates, where this outcome is, on average, equal to her true type. However, she can also withhold part of this “outcome” and show a lower realization than what she actually generates. More formally, $p_M \in \mathcal{M}_t$ if and only if the following two statements are true. First, every $M \in \text{supp}(p_M)$ takes the form $[0, m]$ for some $m \in \mathbf{R}_+$. To state the second property, note that the first implies that any $p_M \in \mathcal{M}_t$ corresponds to a probability distribution over \mathbf{R}_+ where if the realization of this random variable is m , this means the message set is $[0, m]$. The second property is that for any $p_M \in \mathcal{M}_t$, the expectation of this associated random variable is t . That is, in the case where p_M has a finite support,

$$\sum_{[0, m] \in \text{supp}(p_M)} p_M([0, m])m = t.$$

In our example, the agent wants to persuade the principal that her type is large, so it is natural to conjecture that the option of showing a lower outcome will never be used by the agent and hence is irrelevant. In fact, one of our results will be that only the upper bound of a given evidence set will be shown by the agent in an optimal mechanism. However, this result will be entirely independent of the preferences of the agent — the same is true even in a problem where the agent wants to persuade the principal that her type is small.

A special case of the evidence–acquisition model is where the agent has no choice of what message to send at the last step. More formally, this special case is when for every $t \in T$ and every $p_M \in \mathcal{M}_t$, every $M \in \text{supp}(p_M)$ is a singleton. For convenience, we write this special case, the *signal–choice model*, differently. Instead of referring to agent’s choices as evidence acquisition strategies, we refer to them as *signal distributions*. Specifically, the set of options available to type $t \in T$ is a nonempty set $S_t \subseteq \Delta(\mathcal{L})$. We refer to an $s \in \Delta(\mathcal{L})$ as a *signal distribution*. The interpretation is that if the agent chooses $s \in \Delta(\mathcal{L})$, then the principal sees message $m \in \mathcal{L}$ with probability $s(m)$. Equivalently, we can think of this as the singleton message in the realized message set. Letting $S = \cup_t S_t$, with some abuse of notation, we denote the utility functions of the agent and principal by $u : T \times S \times X \rightarrow \mathbf{R}$ and $v : T \times S \times X \rightarrow \mathbf{R}$ respectively. While we will vary the arguments of the utility functions across models, we will always denote the agent’s utility function by u and the principal’s by v .

Similarly to our comments above about the observability of p_M , the model allows the possibility that the realized m reveals the agent’s choice of s always or reveals it with some probability or reveals it for some s choices but not others.

While we discuss the details of games or mechanisms below, the “technology” of evidence/signals imposes some timing structure. Specifically, in both models, we assume

the agent knows her type at the outset. There may be cheap talk between the principal and the agent before the agent chooses a strategy of evidence gathering or a signal distribution. After this, the agent sees the realization of her strategy. In the evidence–acquisition case, this is a set of evidence and (perhaps after further cheap talk) she can then send evidence to the principal. In the signal–choice model, the principal also sees the realization. After this, the principal chooses $x \in X$.

Running Example, Part 3. For a signal–choice version of our running example, we can “convert” the same technology as in the evidence–acquisition model described in Part 2 into a signal–choice model. More precisely, note that the agent in the evidence–acquisition model can pick a distribution over evidence sets and decide what message she will use from each set. In other words, she can have a particular distribution over sets of the form $[0, m]$ and decide for each upper bound m what message $m' \in [0, m]$ she will send to the principal. Recall that the agent of type t can only generate a distribution over sets of the form $[0, m]$ with the property that the expectation of the upper bound m is t . Given this, it is not hard to see that when we convert to signals, this generates the set of signal distributions with expected value less than or equal to t . In other words, for the signal–choice version of our running example, we assume that S_t , the set of signal distributions for type t , is the set of all probability distributions on \mathbf{R}_+ with expected value less than or equal to t . Thus signal distributions are either unbiased or biased “against” the agent. One can think of this as a stylized model where the agent can give the principal one name of a reference for the principal to contact. References cannot be systematically biased in the agent’s favor, but the agent generally cannot predict exactly what a given reference will say.

The usual model of evidence considers games or mechanism design problems where the agent’s set of feasible messages depends on her type. Thus by presenting a message which is only feasible for a certain set of types, the agent proves that her type is in this set. For early contributions in game theory, see Grossman (1981), Milgrom (1981), and Dye (1985). For early contributions in mechanism design theory, see Green and Laffont (1986). For examples of recent papers in game theory or mechanism design, see Shin (2003), Acharya, DeMarzo, and Kremer (2011), Ben-Porath and Lipman (2012), Kartik and Tercieux (2012), Guttman, Kremer, and Skrzypacz (2014). Finally, for closely related work related to both games and mechanisms, see Glazer and Rubinstein (2004, 2006), Hart, Kremer, and Perry (2017), and Ben-Porath, Dekel, and Lipman (2019).

It seems very natural to extend the usual evidence model to incorporate evidence acquisition.

The idea of noisy evidence is also quite natural in some settings. For example, when a lawyer calls a witness to the stand, she may know more about what the witness will say than the judge does, but may not be able to perfectly predict the witness’ testimony.

Similarly, as discussed above, when an agent gives the name of a recommender to the principal, she may not know exactly what the recommender will say. In either case, the evidence that results is a random variable, but does depend on the agent’s true type.

Another interpretation of the signal–choice model is as a natural variation and generalization of the classic Holmstrom (1999) career–concerns model. In the usual career–concerns model, an agent of unknown ability chooses actions which have outcomes that depend on her ability. Both the agent and the principal (or market) learn about the agent’s ability over time by observing these outcomes. Our model is different in that we assume the agent knows her type, while the usual model has symmetric uncertainty about the type.⁴ Our model also generalizes the usual career–concerns model by allowing a richer formulation of actions which allows partial observability of actions as explained above.

The signal–choice model can also be thought of as an “informed agent” version of the Bayesian persuasion model of Kamenica–Gentzkow (2011). As in the Bayesian persuasion model, the agent chooses an “experiment” which reveals information to the principal. Our model differs from Kamenica–Gentzkow in three ways. First, we do not necessarily allow for any possible signal to be created. Second, we assume the agent knows her type, even though she may not know the outcome of the experiment.⁵ Third, we assume the principal does not observe the full experiment as in Kamenica and Gentzkow. Specifically, while we can allow the principal to observe the signal choice of the agent as discussed above, he cannot observe the signals that would have been chosen by other types.

This model also is a natural generalization of models of noisy signaling. For example, in a classic paper, Matthews and Mirman (1983) study a privately–informed firm that chooses an unobservable quantity of output. Because of stochastic demand, this leads to a stochastic price, the realization of which is observed by a potential entrant. Thus the firm effectively chooses a probability distribution over what the rival will see and the expected payoff of the firm depends on it’s type, it’s choice, and the rival’s response.

Deb, Pai, and Said (2018) develop another model which can be thought of as a signal–choice model. A forecaster has private information about the quality of the signals she receives about some random variable. She sees a sequence of signals, announcing a prediction about the random variable after each such observation. After this, the realization of the random variable is observed. The principal updates his beliefs about the quality of her information. To embed this in a signal–choice model, the forecaster’s “message” can be thought of as tuple giving the sequence of forecasts together with the realization of the random variable. A choice of a strategy by the forecaster giving her

⁴See Chen (2015) or Halac and Kremer (2017) for career–concerns models where the agent has private information.

⁵For work on Bayesian persuasion with privately informed agents, see Perez–Richet (2014), Hedlund (2017), and Kosenko (2018).

forecasts as a function of the signals she sees generates a probability distribution over such sequences and hence is a signal choice. Deb, Pai, and Said’s result that the optimal mechanism in this setting does not require commitment by the principal is a special case of our results in Section 5. Their proof restricts attention to deterministic mechanisms; our results show that no such restriction is needed.

Matthews and Postlewaite (1985), Che and Kartik (2009), DeMarzo, Kremer, and Skrzypacz (2019), and Ball and Kattwinkel (2019) give models related to our evidence–acquisition model. In the first three of these papers, an uninformed agent chooses a test which may reveal information about her type. If the test does not produce a result, the agent’s only option is to say there was no result. Otherwise, she can show the result or claim to have none. (In Matthews and Postlewaite, testing always yields a result, but the agent can claim not to have tested.) Thus the agent’s choice of a test, like the choice of an action in our model, produces a probability distribution over a set of options for the agent to reveal. Ball and Kattwinkel take a mechanism design approach to a similar model but where the agent knows her type and where the principal can play a role in determining what test the agent takes.

Given the broad range of models with close connections to the evidence–acquisition and/or signal–choice model, one clear direction for research is to understand these connections better. To what extent can we exploit results known for one of these models in the context of others where the connection has not been noticed? To what extent are we repeatedly seeing examples of a unified phenomena that has not yet been identified?

3 Games

There are many timing assumptions one could consider in modeling the interaction between the agent and principal. We focus on the following sequential game.

First, the agent learns her type. In the evidence–acquisition model, she then chooses $p_M \in \mathcal{M}_t$ and $M \subseteq \mathcal{L}$ is realized. She then chooses $m \in M$. If we consider the signal–choice model instead, the agent simply chooses $s \in S_t$ and the realization m is determined. Either way, the principal observes m but not the agent’s type or other information. That is, in the evidence–acquisition model, the principal does not observe the agent’s choice of p_M or the realization M . In the signal–choice model, the principal does not observe the agent’s choice of s . The principal then chooses $x \in X$.

Given the way we have defined the game, it is straightforward to show that the signal–choice model is a reduced form of the evidence–acquisition model. In the evidence–acquisition model, we can think of the agent choosing p_M and simultaneously choosing

her *messaging strategy* — that is, her strategy for which message m to send as a function of the realization of the message set M . As we vary the agent’s choice of distribution and messaging strategy, we trace out a set of probability distributions over messages m that the principal will observe. Thus each distribution and messaging strategy is equivalent to a signal choice. This is exactly the conversion described in Part 3 of our running example. In light of this, we could analyze the game as an evidence–acquisition model or equivalently replace the set of actions and messaging strategies with the set of induced signal distributions and analyze the game as a signal–choice model.

Running Example, Part 4. We illustrate the game with our running example. Given that the evidence–acquisition model reduces to the signal–choice model, we focus only on the latter. Assume the agent has two equally likely types, h and ℓ where $h > \ell > 0$. For the wage–setting version, we assume $X = \mathbf{R}_+$, $u(t, s, x) = x$, and $v(t, s, x) = -(t - x)^2$. That is, the principal chooses a wage, the agent’s utility is equal to the wage, the principal wishes to set the wage equal to the agent’s true productivity, and the signals are costless. For the hiring version, we assume $X = \{0, 1\}$, $u(t, s, x) = x$, and $v(t, s, x) = x(t - \bar{w})$ where $h > \bar{w} > \ell$. In other words, the agent wants to be hired ($x = 1$), while the principal wants to hire the high type but not the low type.

For either version of the model, the following strategies form a perfect Bayesian equilibrium. Type h chooses the signal distribution which puts probability 1 on h , while ℓ chooses a distribution with probability ℓ/h on h and $1 - (\ell/h)$ on 0. The principal’s belief puts probability 1 on ℓ unless the signal he sees is h . If he sees signal h , his belief puts probability $h/(\ell + h)$ on type h . Either way, he chooses his action accordingly. So in the wage–setting version, he chooses $x = \ell$ if he sees any message other than h and sets $x = (h^2 + \ell^2)/(h + \ell)$ otherwise. In the hiring version, he does not hire if he sees any $m \neq h$ and hires if he sees $m = h$ if

$$\frac{h^2 + \ell^2}{h + \ell} > \bar{w},$$

doesn’t hire if the reverse strict inequality holds, and can choose any probability of hiring otherwise. It is easy to see that, given the principal’s strategy, both types want to maximize the probability on signal h and these signal choices do that. So these strategies form an equilibrium.

If the principal can commit to his reaction to the m he observes, then he can achieve his best possible outcome in the wage–setting version. To be specific, suppose the principal commits to choosing $x = m$ if m is either h or ℓ and to choosing $x = 0$ otherwise. Given any s chosen by the agent, the agent’s expected payoff is less than or equal to the expectation of m since the principal’s choice of x is always weakly below m . Since every $s \in S_t$ has expectation t , this implies that the agent’s payoff must be weakly less than t . Since the agent can obtain a payoff of exactly t by choosing the degenerate s which produces $m = t$ with probability 1, we see that this is an optimal reply for the

agent. Clearly, this enables the principal to choose $x = t$ always, thus achieving his highest possible payoff. It is not hard to show that no equilibrium of the game yields the principal this payoff, so the ability to commit strictly improves the principal’s payoff.

On the other hand, commitment does not help the principal in the hiring version. This is demonstrated in Section 4.3 and generalized in Section 5.

4 Mechanism Design

While we can assume any timing of interest for the case of game theory, for mechanism design, it is more standard to identify a timing structure which allows the principal to obtain the highest possible payoff. Using standard Revelation Principle type arguments, one can show that we can restrict attention to a certain class of direct truth-telling mechanisms. However, these mechanisms are rather complex for the signal-choice model and quite involved for the evidence-acquisition model. Henceforth we use the term *protocol* to refer to the sequence of stages of communication in a mechanism.⁶

For the signal-choice model, we have, in effect, an adverse selection problem (the agent’s private knowledge regarding her type), followed by moral hazard (the agent’s unobserved choice of a signal distribution). Thus a variation on Myerson’s Revelation and Obedience Principle identifies the appropriate protocol.⁷ First, the agent reports a type. Then the principal recommends a signal distribution. Finally, the agent chooses some distribution, the principal observes m , and the principal chooses $x \in X$.

In the evidence-acquisition model, though, the problem is much more complex. In effect, we start with adverse selection (the agent’s type), then have moral hazard (the agent’s choice of a distribution over evidence sets), followed by more adverse selection (the realized set of evidence messages). Consequently, we start similarly to the signal choice case where the agent reports her type, the principal recommends an action, and the agent chooses an action. But after this, the agent makes a report of the realized evidence set, the principal recommends a message choice from this set, and the agent sends a message. Only then does the principal choose $x \in X$. One can show by examples (omitted for brevity) that, in general, each of these steps may be necessary for the principal to obtain the highest possible payoff.

⁶Gerardi and Myerson (2007) have shown that the Revelation Principle may not hold for sequential equilibrium in dynamic environments, raising questions about our multi-stage mechanisms. However, results in Sugaya and Wolitzky (2018)’s Section 4 show that such problems do not arise in our single-agent setting.

⁷For similar results in the evidence literature, see Bull and Watson (2007) and Deneckere and Severinov (2008).

In this section, we give conditions under which we can identify the principal’s recommendations in an optimal mechanism based only on the evidence/signal structure, using little to no information about preferences. Under these conditions, we can eliminate some of the above steps, greatly simplifying the class of mechanisms we need to consider and thus greatly simplifying the analysis.

We begin with the evidence–acquisition model. After stating the protocol formally, we give a condition under which we can identify for each feasible set of messages, a particular message that the principal can always request in an optimal mechanism. Under this condition, we do not need the agent to report the feasible set of messages since the principal’s response to this report is known. Consequently, when this condition holds, we can reduce the evidence–acquisition model to a signal–choice model. After showing this, we develop a condition under which we can identify the principal’s recommended signal distribution for each type report, again largely independently of the preferences. Under this condition, we can then eliminate the principal’s recommendation of a signal distribution, leaving us with a greatly simplified mechanism design problem.

We define the class of deterministic mechanisms for the evidence–acquisition model. However, the principal will typically choose a randomization over these mechanism, for reasons we explain below.

The protocol for evidence–acquisition models has seven stages. We refer to this as the *full protocol for evidence–acquisition models*. Throughout, we use b ’s to denote the agent’s pure strategies at various stages and g ’s to denote the principal’s pure strategies. Recall that \mathcal{M} is the collection of M such that there exists t and $p_M \in \mathcal{M}_t$ with $M \in \text{supp}(p_M)$.

Stage 1. The agent reports a $t \in T$. Denote a reporting strategy for the agent as a function $b_T : T \rightarrow T$ and let B_T denote the set of such functions.

Stage 2. Given the report, the principal requests a distribution p_M over evidence sets. Let $g_M : T \rightarrow \mathcal{M}$, where $g_M(t) \in \mathcal{M}_t$, denote a pure strategy for the principal for this stage and let G_M denote the set of pure strategies.

Stage 3. The agent chooses some feasible p_M and the evidence set M is realized. Let $b_M : T \times T \times \mathcal{M} \rightarrow \mathcal{M}$ denote the agent’s choice of a distribution over evidence sets as a function of her true type, her type report, and the distribution recommended by the principal, where we require $b_M(t, r, p_M) \in \mathcal{M}_t$. Let B_M denote the set of such strategies.

Stage 4. The agent makes a report $\hat{M} \in \mathcal{M}$ of her realized message set. Let this reporting strategy be denoted $b_{\mathcal{M}} : T \times T \times \mathcal{M} \times \mathcal{M} \times \mathcal{M} \rightarrow \mathcal{M}$ where $b_{\mathcal{M}}(t, r, p_M, p'_M, M)$ is the message set the agent reports when her true type is t , her report at Stage 1 was r , the principal requested distribution p_M at Stage 2, she chose distribution p'_M at Stage 3, and the resulting set of feasible messages is M . Let $B_{\mathcal{M}}$ be

the set of such functions.

Stage 5. The principal proposes a message $m \in \hat{M}$ for the agent to send. Let the principal's pure strategy be denoted $g_{\mathcal{L}} : T \times \mathcal{M} \times \mathcal{M} \rightarrow \mathcal{L}$ where $g_{\mathcal{L}}(t, p_M, \hat{M})$ is the message the principal requests when t is the type reported by the agent, p_M the distribution recommended at Stage 2 by the principal, and \hat{M} the message set claimed by the agent at Stage 4. We require $g_{\mathcal{L}}(t, p_M, \hat{M}) \in \hat{M}$. Let $G_{\mathcal{L}}$ denote the set of these pure strategies.

Stage 6. The agent sends a message \hat{m} from the set of available messages. Let this reporting strategy be denoted $b_{\mathcal{L}} : T \times T \times \mathcal{M} \times \mathcal{M} \times \mathcal{M} \times \mathcal{M} \times \mathcal{L} \rightarrow \mathcal{L}$ where $b_{\mathcal{L}}(t, r, p_M, p'_M, M, \hat{M}, m)$ is the message the agent reports when her type is t , her type report is r , the principal requested p_M , she chose p'_M , the resulting message set was M , she reported \hat{M} , and the principal requested message m . Of course, we require that $b_{\mathcal{L}}(t, r, p_M, p'_M, M, \hat{M}, m) \in M$. Let $B_{\mathcal{L}}$ denote the set of such functions.

Stage 7. The principal chooses an outcome as a function of the history he has observed. A pure strategy for the principal is denoted $g_X : T \times \mathcal{M} \times \mathcal{M} \times \mathcal{L} \times \mathcal{L} \rightarrow X$ where $g_X(t, p_M, \hat{M}, m, \hat{m})$ is the principal's choice when the agent reported type t , the principal recommended p_M , the agent reported evidence set \hat{M} , the principal requested message m , and the agent sent message \hat{m} . Let G_X denote the set of such functions.

Let the principal's set of pure mechanisms or pure strategies be denoted $G = G_{\mathcal{M}} \times G_{\mathcal{L}} \times G_X$ with typical element $g = (g_{\mathcal{M}}, g_{\mathcal{L}}, g_X)$. Let $\Gamma = \Delta(G)$ with typical element γ . We let $(\gamma_{\mathcal{M}}, \gamma_{\mathcal{L}}, \gamma_X)$ denote the equivalent behavior strategy to γ . Let $B = B_T \times B_{\mathcal{M}} \times B_{\mathcal{M}} \times B_{\mathcal{L}}$ denote the agent's set of pure strategies with typical element $b = (b_T, b_{\mathcal{M}}, b_{\mathcal{M}}, b_{\mathcal{L}})$. Let $\beta \in \Delta(B)$ denote a typical mixed strategy for the agent.

A version of the standard Revelation Principle for this class of models says that without loss of generality, we can restrict attention to mechanisms where it is optimal for the agent to report truthfully and to obey the principal's recommendations at every stage.

To define incentive compatibility more precisely, note that any (β, γ, t) induces a probability distribution over *complete outcomes* — that is, realized (p_M, x) pairs. We denote this distribution by $\mu(p_M, x \mid \beta, \gamma, t)$. Let $\mathcal{U}(\beta, \gamma, t)$ denote the agent's expected utility in the mechanism γ given strategy β when her type is t or

$$\mathcal{U}(\beta, \gamma, t) = \int_{(p_M, x) \in \mathcal{M} \times X} u(t, p_M, x) d\mu(p_M, x \mid \beta, \gamma, t).$$

We say that a pure strategy $\hat{b} = (\hat{b}_T, \hat{b}_{\mathcal{M}}, \hat{b}_{\mathcal{M}}, \hat{b}_{\mathcal{L}})$ is *truthful and obedient* if for all t, p_M, M , and m , we have $\hat{b}_T(t) = t$, $\hat{b}_{\mathcal{M}}(t, t, p_M) = p_M$, $\hat{b}_{\mathcal{M}}(t, t, p_M, p_M, M) = M$, and

$\hat{b}_{\mathcal{L}}(t, t, p_M, p_M, M, M, m) = m$. That is, the agent reports truthfully and obeys the principal at all stages. Throughout, we use \hat{b}^* to denote any such honest and obedient strategy.⁸

We say that a mechanism γ for the evidence–acquisition model is *incentive compatible* if for every t ,

$$\mathcal{U}(\hat{b}^*, \gamma, t) \geq \mathcal{U}(b, \gamma, t), \quad \forall b \in B$$

for any truthful and obedient strategy \hat{b}^* . (Clearly, this condition also implies that \hat{b}^* is a better strategy for the agent than any mixed strategy $\beta \in \Delta(B)$.)

Given any incentive compatible γ , let $\mu^*(p_M, x \mid \gamma, t) = \mu(p_M, x \mid \hat{b}^*, \gamma, t)$. We refer to μ^* as the *mechanism outcome*.

4.1 Identifying the Recommended Message

Clearly, this is a complex protocol, giving us a complex set of mechanisms and incentive compatibility constraints. In the rest of this section, we introduce two ways to simplify the protocol and conditions under which these simplifications are without loss of generality.

In both cases, the idea is to identify some choices by the principal in a way which depends on the evidence structure but does not depend on the preferences of the principal or the agent. As we will see, the ability to identify such choices greatly reduces the complexity of the protocol and the mechanism design problem.

The idea behind the first simplification is to identify the principal’s response at Stage 5. If for every possible \hat{M} , there is a specific $m \in \hat{M}$ that the principal will always ask for, regardless of the preferences or other details of the model, then we can delete Stage 5, taking as given that the principal will request this message. This enables us to eliminate Stage 4 since the agent’s report of a message set is needed only to give the principal the opportunity to make such a recommendation. Hence we can combine Stages 3 and 6, skipping Stages 4 and 5.

One way to understand when we can identify the principal’s response in this way is by comparison to the literature with exogenously given evidence. In such models, one may need the principal to randomize over which message to request in response to the agent’s type report. The idea is to prevent the agent from knowing how the principal will check various possible lies, thus deterring misreporting. See Glazer and Rubinstein (2004) for illustrative examples. As shown by Bull and Watson (2007), though, under

⁸Note that there are many such strategies since we do not specify how the agent behaves on histories inconsistent with her strategy.

a condition they call normality which Lipman and Seppi (1995) had previously called the full reports condition, this request by the principal is not needed. Normality or full reports says that the agent has available a message which reveals as much information as all the messages the agent has available. Thus asking for this message is the “best” way to deter lies.

We generalize this property to evidence–acquisition models as follows. We say that the evidence technology satisfies *normality* if for every $M \in \mathcal{M}$, there exists $m_M^* \in M$ such that for every $M' \in \mathcal{M}$, we have

$$m_M^* \in M' \iff M \subseteq M'.$$

We refer to the message m_M^* as the *maximal evidence* for M .

To understand this condition, note that $M \subseteq M'$ trivially implies $m_M^* \in M'$ since $m_M^* \in M$. However, we write the condition this way to emphasize the following idea. Intuitively, the only thing that presenting a particular message m proves to the principal is that the agent is able to present this message — that is, that the set of messages the agent has available includes m . In this sense, the presentation of m is evidence directly about M' , the agent’s set of evidence, not about t . It provides evidence only indirectly about t since types differ in terms of which evidence sets they are able or likely to obtain. The definition says that learning that m_M^* is feasible reveals exactly the same information about the agent’s set of messages as learning that every message in M is feasible. In this sense, showing m_M^* reveals exactly what showing every message in M would reveal.

Running Example, Part 5. In our example, \mathcal{M} contains every interval of the form $[0, m]$ for $m \in \mathbf{R}_+$ since each such interval can be generated with positive probability by some (actually either) type. Hence it is easy to see that the most informative message, m_M^* , for the interval $[0, m]$ is the upper bound, m . That is, $m_{[0, m]}^* = m$ or, equivalently, $M = [0, m_M^*]$. This is true as for any $m' \in \mathbf{R}_+$, we have $m_M^* \in [0, m']$ if and only if $[0, m_M^*] \subseteq [0, m']$. Hence our running example satisfies normality. As Theorem 1 below will indicate, this means that only the upper bounds of the intervals will ever be used in an optimal mechanism, regardless of the preferences, as asserted earlier.

To see more concretely that normality is about the information content of messages regarding the set of available messages, consider the following example.

Example 1. The agent has two types, t_1 and t_2 . Each type has only one distribution over evidence sets. Type t_1 obtains evidence set $\{m_1\}$ with probability 1/3, $\{m_2\}$ with probability 1/3, and $\{m_1, m_2\}$ with probability 1/3. Type t_2 receives evidence set $\{m_2\}$ with probability 1. This evidence technology violates normality. First, note that any singleton evidence set trivially has a maximal evidence message since if $M = \{m\}$, then it is obviously true that for any M' , $m \in M'$ iff $M \subseteq M'$. So if normality fails, it is because $\{m_1, m_2\}$ has no maximal evidence message. It is easy to see that this is the

case. For either message $m' \in \{m_1, m_2\}$, the singleton $\{m'\}$ is also an element of \mathcal{M} . Clearly, then, m' cannot be maximal since $m' \in \{m'\}$ but $\{m_1, m_2\} \not\subseteq \{m'\}$.

To see why this is surprising, note that if the agent presents m_1 to the principal, she proves that her type is t_1 as type t_2 never has this message available. Yet m_1 is not maximal evidence from $\{m_1, m_2\}$. Intuitively, presentation of m_1 proves the agent's type but presenting both m_1 and m_2 would prove more about the agent's available messages than m_1 proves.

One way to understand this is to observe that in standard evidence models, the agent's type identifies exactly her set of available messages. In a sense, the agent's *full* type is the pair (t, M) where M is the set of messages the agent has. So proving what t is does not prove the agent's full type.⁹

The following theorem shows that normality will enable us to identify the principal's message recommendations, a result we can then use to simplify the protocol. Recall that a mechanism for the principal is a probability distribution γ over G with associated behavior strategy representation $(\gamma_M, \gamma_L, \gamma_X)$.

Theorem 1. *In the evidence-acquisition model, fix any incentive compatible mechanism γ . If the evidence technology is normal, then there exists an incentive compatible mechanism $(\gamma_M^*, \gamma_L^*, \gamma_X^*)$ with the following properties. First, $\gamma_L^*(t, p_M, M)(m_M^*) = 1$. That is, the principal always recommends the maximal evidence message for any reported M . Second, for all t ,*

$$\mu^*(\cdot \mid \gamma, t) = \mu^*(\cdot \mid \gamma^*, t),$$

so γ and γ^* have the same mechanism outcome for every $t \in T$.

This simplification is, in general, not possible when the evidence technology is not normal. In Appendix B, we analyze Example 1 above with a non-normal evidence technology and show that we cannot identify the message the principal requests from $\{m_1, m_2\}$ independently of the preferences of the principal and agent. More specifically, we give an example of preferences for which it is better for the principal to request m_1 and an example where it is better for him to request m_2 , *even though m_1 perfectly reveals her type*.

Theorem 1 implies that we can simplify the protocol under normality. Since the principal can always recommend the maximal evidence message for any reported message set, we do not need to include the stage where he makes this recommendation. Similarly,

⁹Another way to see this point is to redefine the type space to be the set of possible (t, M) and the set of feasible messages for “type” (t, M) to be M . Applying the standard definition of normality to this model yields our definition.

this means we do not need the agent to report the message set since the mechanism does not depend on it.

Hence a corollary to Theorem 1 is that we can use a simpler protocol. We refer to the following as the *abbreviated protocol for evidence–acquisition models*:

Stage 1. The agent reports a $t \in T$.

Stage 2. Given the report, the principal recommends a distribution over evidence sets for the agent.

Stage 3. The agent chooses a distribution and the evidence set M is realized.

Stage 4. The agent sends a message m from the set of available messages M .

Stage 5. The principal chooses an outcome as a function of the history he has observed, the agent’s report, the recommended distribution, and the message m .

We abuse notation by using the same notation to denote strategies for this protocol. Hence a pure strategy for the agent is now $b = (b_T, b_M, b_L)$ where $b_T : T \rightarrow T$ and $b_M : T \times T \times \mathcal{M} \rightarrow \mathcal{M}$ as before. Also, $b_L : T \times T \times \mathcal{M} \times \mathcal{M} \times \mathcal{M} \rightarrow \mathcal{L}$ where $b_L(t, r, p_M, p'_M, M) \in M$ gives the message the agent sends as a function of her true type t , her reported type r , the principal’s recommended distribution p_M , the distribution she actually chose p'_M , and the realized set M . A pure strategy for the principal is $g = (g_M, g_X)$ where $g_M : T \rightarrow \mathcal{M}$, with $g_M(t) \in \mathcal{M}_t$ and $g_X : T \times \mathcal{M} \times \mathcal{L} \rightarrow X$ gives the principal’s choice of x as a function of the agent’s report, the recommended distribution, and the observed message. Again, we denote the agent’s pure strategies by $B = B_T \times B_M \times B_L$ and the principal’s pure strategies by $G = G_M \times G_X$.

The definition of incentive compatibility for this class of mechanisms is similar to the preceding. Briefly, incentive compatibility requires that an optimal strategy for the agent is to report t truthfully (so $b_T(t) = t$), to follow the principal’s recommendation (so $b_M(t, t, p_M) = p_M$), and to use maximal evidence (so $b_L(t, t, p_M, p_M, M) = m_M^*$).

We have the following corollary, proved in Appendix C:

Corollary 1. *Assume the evidence technology is normal. Then for any incentive compatible mechanism in the full protocol for evidence–acquisition models, there is an incentive compatible mechanism for the abbreviated protocol with the same mechanism outcome.*

Remark 1. The abbreviated protocol simplifies the full protocol by eliminating stages, but also by omitting any recommendation by the principal regarding messaging. One could consider an intermediate protocol where the principal responds to the type report of the agent by recommending both an action and a messaging strategy for what message

the agent should send as a function of her realized message set. It is not hard to give examples where this intermediate form of simplification is possible even when the evidence technology is not normal. For example, this is possible for the mechanism discussed in Appendix B using the assumptions of Example 2. This intermediate simplification also allows the reduction to the signal–choice model discussed next in Section 4.2. ■

4.2 Reduction to Signal Choice

The identification of the principal’s recommended message under normality enables us to reduce the mechanism design problem for the evidence–acquisition model to the mechanism design problem for the signal–choice model. To show this, we first describe the latter. In this case, it is easy to see that we can assume the following *protocol for signal–choice*:

Stage 1. The agent reports a $t \in T$. Let a reporting strategy be denoted $b_T : T \rightarrow T$.

Stage 2. Given the report, the principal requests a signal distribution. A pure strategy is denoted $g_S : T \rightarrow S$.

Stage 3. The agent chooses a signal distribution s and the resulting message is seen by the principal. Let $b_S : T \times T \times S \rightarrow S$ with $b_S(t, r, s) \in S_t$ denote a typical pure strategy for the agent.

Stage 4. The principal chooses an outcome as a function of what has been said. Let $g_X : T \times S \times \mathcal{L} \rightarrow X$ denote a typical pure strategy for the principal at this stage.

Abusing notation, again let $B = B_T \times B_S$ denote the set of pure strategies for the agent and $G = G_S \times G_X$ the set of pure strategies for the principal in this protocol. By the Revelation Principle, we can focus on mechanisms $\gamma \in \Gamma$ with the property that any strategy $\hat{b}^* = (\hat{b}_T^*, \hat{b}_S^*)$ for the agent satisfying $\hat{b}_T^*(t) = t$ and $\hat{b}_S^*(t, t, s) = s$ is a best reply for the agent to γ . Again, we refer to any such \hat{b}^* as truthful and obedient. Given an incentive compatible mechanism γ , we can define the mechanism outcome as the function mapping t to probability distributions over outcomes, here defined as (s, x) pairs. I.e., we can write $\mu^*(s, x \mid \gamma, t)$ as the probability distribution over (s, x) induced by the strategies (\hat{b}^*, γ) given the agent’s type is t .

Just as in our analysis of games in Section 3, we can think of the agent’s strategy in the evidence–acquisition model as a choice of a distribution over evidence sets and a messaging strategy. Again, any given distribution and messaging strategy generates a probability distribution over the message the agent shows the principal. Thus we can replace the selection of a distribution and a messaging strategy with the selection of a

signal distribution. In general, this change takes away some of the principal's ability to influence the agent's decisions and will lead to a less effective mechanism. However, under normality, the ability to reduce to the abbreviated protocol implies that this change does not harm the principal.

More formally, fix an evidence-acquisition model. We construct a signal-choice model from it as follows. For any $p_M \in \mathcal{M}$ and any function $\sigma : \text{supp}(p_M) \rightarrow \mathcal{L}$ such that $\sigma(M) \in M$, we can define a signal $s \in \Delta(\mathcal{L})$ by

$$s(m) = p_M(\{M \mid \sigma(M) = m\}).$$

Let $\Sigma(p_M)$ denote the set of such σ functions given p_M and let $s_{(p_M, \sigma)}$ denote the distribution on \mathcal{L} induced by (p_M, σ) . Let

$$S_t = \{s_{(p_M, \sigma)} \mid p_M \in \mathcal{M}_t, \sigma \in \Sigma(p_M)\}.$$

We can define utility functions for the signal-choice model by letting $u(t, s_{(p_M, \sigma)}, x) = u(t, p_M, x)$ and analogously for v . This is exactly the translation from evidence acquisition to signal choice discussed less formally in Section 3.

The following result explains the sense in which the signal-choice model so constructed is equivalent to the evidence-acquisition model under normality.

Theorem 2. *In the evidence-acquisition model, fix any incentive compatible mechanism γ . If the evidence technology is normal, then there exists an equivalent incentive compatible mechanism γ^* in the signal-choice model constructed from it in the following sense. For any truthful and obedient strategy \hat{b}^* for the agent in the signal-choice model given γ^* , we have*

$$\mu^*(p_M, x \mid \gamma, t) = \mu(s_{(p_M, \sigma^*)}, x \mid \hat{b}^*, \gamma^*, t),$$

for all t where σ^* is the function $\sigma^*(M) = m_M^*$ for all $M \in \text{supp}(p_M)$. In this sense, γ and γ^* have equivalent mechanism outcomes for every $t \in T$.

In short, under the assumption of normality, any outcome that can be induced by a mechanism for the evidence-acquisition model can be induced by an incentive compatible mechanism in the protocol for the signal-choice model. This is analogous to our result on games in Section 3.

One can consider mechanisms with a variety of other timings. For example, perhaps the agent only comes to the principal *after* having generated evidence. Recognizing this, the optimal mechanism takes into account the way the rules of the mechanism affect these incentives. For example, this seems like a natural way to think about courts. The lawyers know the rules of the court in advance and work to obtain evidence before bringing the case to court. It is easy to show the analog of Theorem 1, Corollary 1,

and Theorem 2 for this model. More specifically, it is still true that under normality, one can restrict attention to mechanisms for which the principal always recommends the maximal evidence message for any evidence set, enabling us to use (an appropriately modified version of) the abbreviated protocol and reduce to a version of the signal–choice model.

4.3 Identifying the Recommended Signal

In this section, we focus on the signal–choice model, where this can be interpreted as a reduced form of the evidence–acquisition model under normality.

While normality greatly simplifies the mechanism design problem, the problem is still complex. We next turn to conditions under which we can identify the signal choice the principal requests as a function of the type.

For convenience, in this section, we assume \mathcal{L} is finite and write a signal distribution $s \in S$ as a (row) vector in $\mathbf{R}_+^{\#\mathcal{L}}$. Fix t^* and $s^*, \hat{s}^* \in S_{t^*}$. We say that s^* is *more informative than* \hat{s}^* if there exists an $\#\mathcal{L} \times \#\mathcal{L}$ Markov matrix Λ such that $s^* \Lambda = \hat{s}^*$ and for every t and every $s \in S_t$, $s \Lambda \in \text{conv}(S_t)$.¹⁰

In the case where each S_t is finite, we can give an equivalent statement which will aid in clarifying the intuition of this condition. Let \mathcal{S} denote the matrix formed by “stacking” the signal distributions. In other words, this is a matrix with $\#\mathcal{L}$ columns and a number of rows equal to $\sum_t \#S_t$. The first $\#S_{t_1}$ rows are the signal distributions available to t_1 , the next $\#S_{t_2}$ rows those available to t_2 , etc. Note that if $s \in S_t \cap S_{t'}$ for $t \neq t'$, then s appears (at least) twice in the matrix. Then s^* is more informative than \hat{s}^* if there exists a Markov matrix Λ such that $\mathcal{S} \Lambda = \hat{\mathcal{S}}$ where the matrix $\hat{\mathcal{S}}$ has \hat{s}^* in the row corresponding to s^* in \mathcal{S} and for any row s of $\hat{\mathcal{S}}$ corresponding to one of type t ’s signal distributions, we have $s \in \text{conv}(S_t)$.

To see the intuition, recall Blackwell–Girshick’s (1954) (BG) comparison of experiments. In their model, there are n states of the world. An experiment gives a probability distribution over a finite set of observations as a function of the state of the world. If there are N possible observations, we can write this as an $n \times N$ matrix E where e_{ij} is the probability of observation j in state i . Suppose we have two experiments, E and F . BG say experiment E is more informative than experiment F if there exists a Markov matrix Λ such that $E \Lambda = F$. The matrix Λ defines a garbling of the results of experiment E , so this says that F can be obtained from E by adding random noise.

¹⁰A matrix is Markov if all entries are non–negative and every row sum is 1.

Thus we can interpret our informativeness comparison as saying that the “experiment” \mathcal{S} is more informative than “experiment” $\hat{\mathcal{S}}$ in the sense that we can obtain the latter by adding noise to the former. To understand the sense in which \mathcal{S} and $\hat{\mathcal{S}}$ can be thought of as experiments, note that the rows in an experiment correspond to states of the world, while a row in \mathcal{S} corresponds to a (type, signal distribution) pair. Intuitively, just as we can think of (t, M) as the (partly endogenous) “full type” in the evidence–acquisition model, it is natural to think of (t, s) as the “full type” in the signal–choice model.

To see the sense in which the existence of Λ implies s is more informative than s' , suppose we have a mechanism in which the principal recommends s' if the agent reports that her type is t . Suppose the principal changes the mechanism to recommend s in this situation instead and changes no other recommendations. Suppose that the principal’s response to messages he subsequently receives from the agent after this recommendation is to “garble” them according to the Markov matrix Λ and then to respond the way the original mechanism specified. If the agent uses signal s , then the resulting distribution over the garbled message will be $s\Lambda$. By hypothesis, this is s' . Thus the distribution over the principal’s choice of x will be the same as in the original mechanism. Suppose that the agent’s true type is \hat{t} and that she uses some signal $\hat{s} \in S_{\hat{t}}$. Then the induced distribution over garbled messages will be $\hat{s}\Lambda$. By hypothesis, this is an element of $\text{conv}(S_{\hat{t}})$. In other words, in the original mechanism, type \hat{t} could have generated this distribution over messages by a particular randomization over her available signals. Thus the expected outcome this type would generate is something she could have generated in the original mechanism. If the original mechanism was incentive compatible, then this deviation is not profitable. Thus the new mechanism is incentive compatible and generates the same outcome as the original one.

Of course, if the more informative signal is so costly that the agent can never be induced to choose it, then it is not useful. Hence to identify the signal the principal will ask the agent to send, we need a condition which identifies that signal as both informative and not excessively expensive. For simplicity, we eliminate the second issue by assuming signals are costless in the sense that neither u nor v depend on the agent’s choice of s . To be precise, we say the model has *costless signals* if for every $t \in T$, $x \in X$, and $s, s' \in S_t$, we have $u(t, s, x) = u(t, s', x)$ and $v(t, s, x) = v(t, s', x)$.

Theorem 3. *In the signal–choice model with costless signals, fix any incentive compatible mechanism γ with marginal γ_S on G_S . If there exists t^* and $s^*, \hat{s}^* \in S_{t^*}$ such that s^* is more informative than \hat{s}^* , then there exists an incentive compatible mechanism (γ_S^*, γ_X^*) satisfying the following two properties. First,*

$$\gamma_S^*(t)(s) = \begin{cases} \gamma_S(t)(s), & \text{if } t \neq t^* \text{ or } s \notin \{s^*, \hat{s}^*\}; \\ \gamma_S(t^*)(s^*) + \gamma_S(t^*)(\hat{s}^*), & \text{if } t = t^* \text{ and } s = s^*; \\ 0, & \text{if } t = t^* \text{ and } s = \hat{s}^*. \end{cases}$$

That is, γ^* moves any probability on recommending \hat{s}^* for t^* to recommending s^* instead. Second, for all t ,

$$\mu^*(\cdot | \gamma, t) = \mu^*(\cdot | \gamma^*, t),$$

so γ and γ^* have the same mechanism outcome for every $t \in T$.

Theorem 3 implies that in any model with costless signals, if type t has some signal distribution $s^* \in S_t$ which is more informative than any other $s \in S_t$, then the principal may as well always recommend s^* to t . If every t has such a most informative signal distribution, then Stage 2 of the mechanism protocol is not needed as we can restrict attention to mechanisms where every type of the agent is induced to choose her most informative signal distribution. In such a case, we can focus on the following *succinct protocol*:

Stage 1. The agent reports a $t \in T$ and chooses a signal distribution s . Let a reporting strategy be denoted $b_T : T \rightarrow T$ and a signal distribution strategy be $b_S : T \rightarrow S$ with $b(t) \in S_t$.

Stage 2. The principal observes the report, the realized m , and chooses an outcome. Let $g_X : T \times \mathcal{L} \rightarrow X$ denote a typical pure strategy for the principal.

Abusing notation yet again, let $B = B_T \times B_S$ denote the set of pure strategies for the agent and G the set of pure strategies for the principal in this protocol. When each type t has a most informative signal distribution s_t^* , we can focus on mechanisms $\gamma \in \Gamma$ with the property that the strategy $\hat{b}_T(t) = t$ and $\hat{b}_S(t) = s_t^*$ is a best reply for the agent to γ .

Running Example, Part 6. We showed in Part 5 of the example that the evidence–acquisition technology is normal. In particular, given any realized message set of the form $[0, m]$, the upper bound m is the most informative message for the set. Hence Theorem 2 implies that we can focus on the signal–choice model where for each t , S_t is the set of all distributions on \mathbf{R}_+ with expectation less than or equal to t . Since \mathbf{R}_+ is not finite, we need to adjust the example to apply our condition. So let \mathcal{L} be any finite subset of \mathbf{R}_+ containing at least T , where we generalize the example, now letting T be any finite subset of \mathbf{R}_+ , not necessarily $\{\ell, h\}$. Assume S_t is the set of all probability distributions on \mathcal{L} with expectation less than or equal to t .

We now show that the most informative signal distribution for type t is the degenerate distribution on t . Fix any type t^* . Let $s^* \in S_{t^*}$ denote the degenerate distribution putting probability 1 on $m = t^*$ and fix any other $s \in S_{t^*}$. Let the Λ matrix be an identity matrix

but with the row corresponding to $m = t^*$ replaced by s . That is, we let

$$\Lambda = \begin{pmatrix} 1 & 0 & 0 & \dots & 0 & 0 \\ 0 & 1 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \dots & \vdots & \vdots \\ s(m_1) & s(m_2) & s(m_3) & \dots & s(m_{\#\mathcal{L}-1}) & s(m_{\#\mathcal{L}-1}) \\ \vdots & \vdots & \vdots & \dots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 1 & 0 \\ 0 & 0 & 0 & \dots & 0 & 1 \end{pmatrix}.$$

Then $s^*\Lambda = s$. Fix any other type t and any $\hat{s} \in S_t$. Let $\tilde{s} = \hat{s}\Lambda$. For $m \neq t^*$, we have $\tilde{s}(m) = \hat{s}(m) + \hat{s}(t^*)s(m)$. For $m = t^*$, we have $\tilde{s}(t^*) = \hat{s}(t^*)s(t^*)$. So

$$\begin{aligned} \sum_m \tilde{s}(m)m &= \sum_{m \neq t^*} [\hat{s}(m) + \hat{s}(t^*)s(m)]m + \hat{s}(t^*)s(t^*)t^* \\ &= \sum_{m \neq t^*} \hat{s}(m)m + \sum_{m \neq t^*} \hat{s}(t^*)s(m)m + \hat{s}(t^*)s(t^*)t^* \\ &= \sum_{m \neq t^*} \hat{s}(m)m + \hat{s}(t^*) \sum_m s(m)m \\ &\leq \sum_{m \neq t^*} \hat{s}(m)m + \hat{s}(t^*)t^* \\ &= \sum_m \hat{s}(m)m \leq t. \end{aligned}$$

The next-to-last line follows from $s \in S_{t^*}$ and therefore $\sum_m s(m)m \leq t^*$. The last inequality on the last line follows from $\hat{s} \in S_t$ and therefore $\sum_m \hat{s}(m)m \leq t$. Hence for every $\hat{s} \in S_t$, we see that $\hat{s}\Lambda$ is a probability distribution over \mathcal{L} with expectation weakly less than t and hence is an element of S_t and therefore of $\text{conv}(S_t)$. Hence s^* is more informative than s .

Now that we have identified the signal choices for each type in the optimal mechanism, it is not difficult to compute the rest of the mechanism. We already showed that the principal can achieve his best possible outcome for each type when his utility function is $-(t-x)^2$, so consider the ‘‘fixed wage’’ case where the principal’s choice is to hire the agent ($x = 1$) or not ($x = 0$) and his payoff is $x(t-w)$ where $w \in (\ell, h)$. Recall that types are equally likely. The agent’s payoff is x . Let $\gamma^*(t)$ denote the probability the principal chooses $x = 1$ when the agent reports type t and the realized message m also equals t . Given that the mechanism will induce truthful reporting and will induce the agent to choose the degenerate distribution with $m = t$, the principal’s expected payoff is

$$\frac{1}{2} \gamma^*(h)(h-w) + \frac{1}{2} \gamma^*(\ell)(\ell-w).$$

To make it easier to induce the agent to choose the appropriate degenerate distribution, the optimal mechanism has $x = 0$ if the message observed differs from the agent's type report. As we will see, type h never wishes to imitate ℓ , so we do not need to impose this incentive compatibility constraint. Hence the only incentive compatibility constraint we require is

$$\gamma^*(\ell) \geq \gamma^*(h) \frac{\ell}{h},$$

since the maximum probability ℓ can put on $m = h$ is when she chooses the distribution with probability ℓ/h on h and the remaining probability on 0. Since the principal's utility is decreasing in $\gamma^*(\ell)$ and increasing in $\gamma^*(h)$, the constraint is binding. Hence the principal chooses $\gamma^*(h)$ to maximize

$$\gamma^*(h) \left[\frac{1}{2} (h - w) + \frac{1}{2} \frac{\ell}{h} (\ell - w) \right].$$

So if

$$\frac{h^2 + \ell^2}{h + \ell} > w,$$

the optimal mechanism has $\gamma^*(h) = 1$. If we have the opposite strict inequality, it has $\gamma^*(h) = \gamma^*(\ell) = 0$. In both cases, type h has no incentive to imitate type ℓ , as asserted.

Also, in both cases, the outcome is the same as in the equilibrium we computed for this example in Section 3. In this sense, there is no value to the principal from commitment: he obtains the same outcome whether he is able to commit to his responses to the agent as in a particular equilibrium of the game where he cannot commit. We present a generalization of the result of this example in the following section.

5 Commitment

We saw in Section 4 that there is no value to commitment in the hiring version of our running example. In this section, we generalize this result. More specifically, we introduce an assumption on the relationship between the preferences of the principal and of the agent and show that if evidence acquisition is costless, then this assumption implies there is no value to commitment. Thus, under our assumption on the preferences we generalize previous results for deterministic evidence models to a general structure of stochastic evidence.

More specifically, the result does not rely on assumptions regarding normality of the evidence technology or informativeness of signals. It also does not depend on the specific protocol used and so can allow for departures from the protocols discussed in the previous sections. The primary requirement for the protocol is that we compare “apples

to apples” — that is, that we compare what the principal can obtain with commitment in a particular protocol to what he can obtain without commitment in the same protocol. Our primary substantive assumption on the protocol is finiteness, though, as we explain, this can be substantially relaxed. To tractably analyze sequential rationality, we impose further structure on the protocol, but this structure is for convenience, not because the result hinges on it.

In particular, our assumptions on preferences are weaker in some ways, stronger in other ways than the assumptions used in the previous work on value of commitment in games with evidence — see, in particular, Glazer and Rubinstein (2004, 2006), Sher (2011), and Hart, Kremer, and Perry (2017). Unlike these models, we do not require type-independent utility functions. After stating our result, we compare our assumptions to theirs in more detail.

We present our result in two parts. First, we show that under only a finiteness assumption on the protocol, there is a Nash equilibrium with the same outcome as the optimal mechanism whenever the preferences satisfy a certain condition. Second, we show that we can strengthen the conclusion to perfect Bayesian equilibrium under stronger conditions.

To state the first result, fix any protocol for the evidence-acquisition model (including the special case of the signal-choice model). Let B denote the set of pure strategies for the agent and G the set of pure strategies for the principal as before. As before, let $U(\beta, \gamma, t)$ denote the agent’s expected payoff under the protocol given mixed strategy profile (β, γ) when her type is t . Let $U(\beta, \gamma) = E_t U(\beta, \gamma, t)$. Let the principal’s expected utility given that the agent is type t be denoted $V(\beta, \gamma, t)$ and let $V(\beta, \gamma) = E_t V(\beta, \gamma, t)$. Given any $\gamma \in \Delta(G)$, let $BR(\gamma)$ denote the agent’s set of best replies — i.e.,

$$BR(\gamma) = \{\beta \in \Delta(B) \mid U(\beta, \gamma) \geq U(\beta', \gamma), \forall \beta' \in \Delta(B)\}.$$

Let

$$V^* = \max_{\gamma \in \Gamma} \max_{\beta \in BR(\gamma)} V(\beta, \gamma).$$

In other words, V^* is the principal’s maximal expected payoff when he can commit to any mixed strategy in the protocol *and* can choose the agent’s best reply to his strategy. If (β^*, γ^*) solves $V^* = V(\beta^*, \gamma^*)$ and $\beta^* \in BR(\gamma^*)$, we say (β^*, γ^*) is *optimal for the principal*.

For our Nash equilibrium result, the only assumption we make on the protocol and evidence structure is that B and G are finite. As the proof will make clear, it would not be difficult to weaken this finiteness assumption further.

We make two assumptions on preferences. First, we assume that evidence acquisition

is costless. That is, we can write the utility function of the agent and principal as $u : T \times X \rightarrow \mathbf{R}$ and $v : T \times X \rightarrow \mathbf{R}$ respectively, dropping the p_M argument used in Section 2.

Our primary preference condition is that there is some function $\nu : T \rightarrow \mathbf{R}$ such that $v(t, x) = \nu(t)u(t, x)$ for all $(t, x) \in T \times X$. The key implication this has which we will use in the proofs is that that $V(g, b, t) = \nu(t)U(g, b, t)$ for all $(g, b, t) \in G \times B \times T$. When this holds, we say the preferences are *semi-aligned*.

The critical implication of semi-aligned preferences is that they imply that if $U(g, b, t) = U(g', b', t)$ for all $t \in T$, then $V(g, b) = V(g', b')$. That is, if all types of the agent are indifferent between two outcomes, then the principal is as well.

The following is our result on value of commitment relative to Nash equilibrium.

Theorem 4. *Fix any protocol for which B and G are finite. Assume evidence acquisition is costless and that preferences are semi-aligned. If (β^*, γ^*) is optimal for the principal, then there exists $\hat{\beta} \in \Delta(B)$ such that $(\hat{\beta}, \gamma^*)$ is a Nash equilibrium of the game induced by the protocol and $V(\hat{\beta}, \gamma^*) = V^*$.*

Proof. Consider the restricted game where the principal's set of pure strategies is G , but the agent's set of pure strategies is $B \cap BR(\gamma^*)$.

By finiteness of B and G , the restricted game has a mixed equilibrium, say, $(\hat{\beta}, \hat{\gamma})$. Since $\hat{\beta}$ puts all probability on $BR(\gamma^*)$, we know that

$$U(\hat{\beta}, \gamma^*, t) = U(\beta^*, \gamma^*, t)$$

for all t . Hence

$$\begin{aligned} V(\hat{\beta}, \gamma^*) &= \mathbf{E}_t[V(\hat{\beta}, \gamma^*, t)] \\ &= \mathbf{E}_t[\nu(t)U(\hat{\beta}, \gamma^*, t)] \\ &= \mathbf{E}_t[\nu(t)U(\beta^*, \gamma^*, t)] \\ &= V(\beta^*, \gamma^*). \end{aligned}$$

Since $(\hat{\beta}, \hat{\gamma})$ is a Nash equilibrium of the restricted game, we also have $V(\hat{\beta}, \hat{\gamma}) \geq V(\hat{\beta}, \gamma^*)$. We now show that this weak inequality must be an equality.

Let

$$\bar{\gamma}_\varepsilon = \varepsilon\hat{\gamma} + (1 - \varepsilon)\gamma^*.$$

By finiteness of B , if we choose ε sufficiently small, then any best reply to $\bar{\gamma}_\varepsilon$ is necessarily a best reply to γ^* . That is, there exists $\bar{\varepsilon} > 0$ such that for all $\varepsilon \in (0, \bar{\varepsilon})$,

$$BR(\bar{\gamma}_\varepsilon) = \{\beta \in BR(\gamma^*) \mid U(\beta, \hat{\gamma}) \geq U(\beta', \hat{\gamma}), \forall \beta' \in BR(\gamma^*)\}.$$

But by construction, $\hat{\beta}$ is a best reply to $\hat{\gamma}$ when the agent is restricted to using best replies to γ^* . Hence for $\varepsilon \in (0, \bar{\varepsilon})$, $\hat{\beta} \in BR(\bar{\gamma}_\varepsilon)$.

Hence for $\varepsilon \in (0, \bar{\varepsilon})$,

$$V^* \geq V(\hat{\beta}, \bar{\gamma}_\varepsilon) = \varepsilon V(\hat{\beta}, \hat{\gamma}) + (1 - \varepsilon)V(\hat{\beta}, \gamma^*).$$

As shown above, $V(\hat{\beta}, \hat{\gamma}) \geq V(\hat{\beta}, \gamma^*) = V(\beta^*, \gamma^*)$. This implies

$$V^* \geq V(\hat{\beta}, \gamma^*) = V(\beta^*, \gamma^*).$$

But $V^* = V(\beta^*, \gamma^*)$, so the inequality must be an equality, implying $V(\hat{\beta}, \hat{\gamma}) = V(\hat{\beta}, \gamma^*)$.

This implies $(\hat{\beta}, \gamma^*)$ is a Nash equilibrium of the overall (unrestricted) game since it says that γ^* is a best reply to $\hat{\beta}$ (since it gives the principal the same payoff as his equilibrium strategy in the restricted game and he was not restricted in that game) and $\hat{\beta}$ is necessarily a randomization over best replies to γ^* . From the above, $V^* = V(\hat{\beta}, \gamma^*)$. ■

Remark 2. The proof of this result does not use the assumption that evidence acquisition is costless, though the PBE result below does. On the other hand, the assumption that preferences are semi-aligned can be unnatural when evidence acquisition is costly. When preferences are semi-aligned and evidence acquisition is costly, the principal cares about the agent's costs. In some situations, this seems natural. For example, if the principal is a social welfare agency which requires potential aid recipients to provide documentation establishing their need, then the principal may well be concerned about not overburdening these recipients. On the other hand, if there are recipient types to whom the principal does not want to give aid, then we must model this by assuming $\nu(t) < 0$ for these types. In this case, our preference assumption would imply that the principal is made better off when such types bear high costs getting documents, an odd assumption.

To extend the result to perfect Bayesian equilibrium, we need to put more structure on the protocol. Otherwise, it is difficult to characterize what kind of choices the principal might have at certain information sets and therefore difficult to characterize sequential rationality at all information sets.

To economize on notation and keep the proofs as simple as possible, we assume the protocol is a multi-stage game with certain properties. It will be apparent from the proof that our argument does not require as much structure as we impose. To avoid repetition, we state the definitions, result, and proof only for the signal-choice version of the model, but it is straightforward to rewrite all that follows for the evidence-acquisition model.

More specifically, we assume the protocol has K stages. The first K_1 stages involve only exchange of cheap-talk messages. Then the agent chooses a signal distribution where the principal observes the realized message. Then there are another K_2 stages

involving exchange of cheap-talk messages. Finally, the principal chooses x . To be more specific, at each odd-numbered stage $k = 1, 3, \dots, K_1 - 1$, the agent can send a cheap-talk message. At each even-numbered stage $k = 2, 4, \dots, K_1$, the principal can send a cheap-talk message. The set of cheap talk messages available to the speaker at stage k is C_k , where this does not depend on the agent's type or previously sent cheap-talk messages. At each of these stages, the opponent observes the cheap-talk message chosen by the speaker. At stage $K_1 + 1$, the agent chooses a signal $s \in S$, where the agent of type t is restricted to choosing $s \in S_t$. The principal does not observe the signal distribution, but both the principal and the agent observe its realization. At each odd-numbered stage $k = K_1 + 3, K_1 + 5, \dots, K_1 + K_2 + 1$, the agent chooses a cheap-talk message. At each even-numbered stage $k = K_1 + 2, K_1 + 4, \dots, K_1 + K_2$, the principal chooses a cheap-talk message. Again, the set of cheap talk messages available to the speaker at stage k is C_k , where this does not depend on the agent's type, the agent's signal choice at stage $K_1 + 1$ or its realization, or previously sent cheap-talk messages. Again, at each of these stages, the opponent observes the cheap-talk message chosen by the speaker. Finally, at stage $K_1 + K_2 + 2$, the principal chooses $x \in X$.

We say that a protocol satisfying these properties is *simple*. We say a simple protocol is *very simple* if $K_2 = 0$ — that is, if the principal chooses $x \in X$ immediately after observing the realization of the agent's chosen signal.

Theorem 5. *Given any very simple protocol, under the assumptions of Theorem 4, then there is a perfect Bayesian equilibrium $(\hat{\beta}, \hat{\gamma})$ with $V(\hat{\beta}, \hat{\gamma}) = V^*$.*

We conclude this section by discussing the strength and tightness of our assumptions. First, we compare them to those used in earlier results showing no value to commitment in mechanism design problems with evidence. As we noted at the outset, the previous literature all considered deterministic evidence, so the main way our result differs is in extending to a general model of stochastic evidence. The comparison of the assumptions on preferences is more involved.

The first results in the deterministic evidence literature were shown by Glazer and Rubinstein (2004, 2006). They considered problems where the principal chooses between two outcomes, called accept ($x = 1$) and reject ($x = 0$). The agent's utility is x . The principal's utility is x if the agent's type is in a certain set of types, $-x$ otherwise. They consider nonstochastic evidence — that is, each type has only a single degenerate distribution over evidence sets. They do not assume the deterministic evidence version of normality, just as we do not require normality, but, unlike us, they assume a particular protocol. As for our assumptions on preferences, if we assume $\#X = 2$ and that there are no costs to signals/distributions, then our assumption of semi-aligned preferences is without (further) loss of generality. To see this, note that in the costless case, we can write the agent's utility function as $u(t, x)$. When there are only two x 's, say, x_0 and x_1 , we can always renormalize the agent's payoffs so that $u(t, x_0) = 0$ for all t , $u(t, x_1) = 1$ for

types t who prefer x_1 to x_0 , and $u(t, x) = -1$ for types who prefer x_0 to x_1 .¹¹ We can also renormalize the principal’s utility function so that $v(t, x_0) = 0$ for all t . Without loss of generality, assume $v(t, x_0) \neq v(t, x_1)$ for all t .¹² Hence, given these renormalizations, we can write $v(t, x) = \nu(t)u(t, x)$ by defining $\nu(t) = v(t, x_1)/u(t, x_1)$ for all t . This preference structure is more general than that of Glazer–Rubinstein since we are not imposing their assumption that u is independent of t .

Sher (2011) generalizes Glazer–Rubinstein to allow a finite set of actions but where the principal’s utility can be written as a concave function of the agent’s utility. Hart, Kremer, and Perry (2017) generalize the preference conditions by assuming that the principal’s utility can be written as a single-peaked function of the agent’s utility given any belief over the agent’s type. Both papers continue to assume that the agent’s utility is independent of her type. Hart, Kremer, and Perry restrict the set of mechanisms to deterministic mechanisms. They also impose the version of normality for deterministic evidence models, unlike Sher and unlike us. Our assumptions on the protocol are weaker than theirs and our assumptions on preferences do not impose type-independence. On the other hand, if we specialize our assumptions to the type-independent case, then their assumptions on preferences are weaker than ours. Thus neither model is strictly more general than the other.

Our preference assumptions are quite strong in the case where the agent’s utility is type independent. In this case, we may as well assume that there are only two outcomes. In other words, when the agent’s utility is simply $u(x)$ and the principal’s utility is $\nu(t)u(x)$, we may as well assume that there are only two choices of x available to the principal. Intuitively, if the principal chooses a mechanism which pools a certain set of types together, the principal’s payoff given this pool will depend on the expectation of $\nu(t)$ across the pool. If this is positive, then the best x for the principal is the one which maximizes $u(x)$, while if it is negative, the best x minimizes $u(x)$. Hence the principal will never strictly prefer any action to one of these two. We emphasize, though, that we do *not* require type-independent utility for the agent, so other actions may be relevant to the principal.

Finally, we discuss the tightness of our assumptions. We have shown no value to commitment when we can write the principal’s payoff as a function of t times the agent’s payoff. It’s immediate that the same holds if we can write the principal’s payoff as a function of t times the agent’s payoff plus some function of t . That is, writing this for the signal-choice model for simplicity, the result still holds if we can write $v(t, x) = \nu(t)u(t, x) + \varphi(t)$. This is immediate since the function $\varphi(t)$ does not affect any optimal

¹¹Types who are indifferent between x_0 and x_1 do not affect the arguments, so we can assume without loss of generality that there are no such types.

¹²If this is violated for some t , then the principal’s decisions are the same as those he would make if such t were impossible. Hence such types can be disregarded.

choices. Put differently, we could renormalize the principal’s payoff by subtracting the expectation of this function.

On the other hand, the result does not generalize to assuming $v(t, x) = \nu(t)u(t, x) + \varphi(t) + \psi(x)$. This fact is demonstrated by Part 4 of our running example. Recall that we considered a version where $u(t, x) = x$ and $v(t, x) = -(x - t)^2$ and showed that commitment enabled the principal to obtain a strictly higher payoff than in any equilibrium. Note that the principal’s payoff function can be written as $2tx - t^2 - x^2 = 2tu(s, x, t) - t^2 - x^2$. So letting $\nu(t) = 2t$, $\varphi(t) = -t^2$, and $\psi(x) = -x^2$, we see that this shows the result does not extend. By contrast, the no-value-to-commitment result in Ben-Porath, Dekel, and Lipman (2019) covers this example in the case of deterministic evidence.

Similarly, the result does not extend to multiple agents. For example, suppose we have two agents, $i = 1, 2$. Suppose the principal’s decision is which agent to give one unit of a good to. Let $X = \{0, 1, 2\}$ where $x = 0$ means the principal keeps the good and $x = i$ means the principal gives the good to agent i . Suppose agent i ’s utility function, $u_i(t_i, x)$ is 1 if i receives the good, 0 otherwise. Suppose the principal’s payoff is $v_i(t_i)$ if he gives the good to agent i . Then we can write the principal’s utility as $v(t, x) = \sum_i v_i(t_i)u_i(t_i, x)$, a natural generalization of our assumption of semi-aligned preferences for the multiple agent case. However, in Appendix G, we give an example showing that the no-value-to-commitment result does not hold for this model even though the principal has only two actions.¹³ Again, this is in contrast to results in Ben-Porath, Dekel, and Lipman (2019) for deterministic evidence.

6 Conclusion

Our results are primarily methodological. While our running example suggests how these results may be useful in practice, an important next step is to apply these tools to develop new insights, e.g., about incentives for evidence acquisition. While there are some results on how optimal mechanisms may introduce distortions to provide incentives for *information* acquisition, there has been very little analysis of *evidence* acquisition.

As discussed above, there is a very broad range of models with close connections to the evidence-acquisition and/or signal-choice model. In light of this, one clear direction for research is to understand these connections better with the goal of better unifying our understanding of these issues.

¹³One can show that there is still a value to commitment if we allow the principal to keep the good.

Appendix

A Proof of Theorem 1

Fix any incentive compatible mechanism $(\gamma_{\mathcal{M}}, \gamma_{\mathcal{L}}, \gamma_X)$. We show how to construct an incentive compatible mechanism with the same mechanism outcome with the property that the principal always recommends m_M^* when the agent reports message set M .

Fix any profile $(\hat{t}, \hat{p}_M, \hat{M}, \hat{m})$ consisting of a type report $\hat{t} \in T$, a recommended distribution over evidence sets $\hat{p}_M \in \text{supp}(\gamma_{\mathcal{M}}(\hat{t}))$, a reported message set $\hat{M} \in \mathcal{M}$, and a requested message $\hat{m} \in \text{supp}(\gamma_{\mathcal{L}}(\hat{t}, \hat{p}_M, \hat{M}))$ such that $\hat{m} \neq m_M^*$. If there is no such tuple, then the principal always recommends maximal evidence, so there is nothing to prove. We construct an alternative mechanism which replaces the recommendation \hat{m} with a recommendation of m_M^* in this situation and will show that this mechanism is incentive compatible and implements the same outcome as the original mechanism. For brevity, let $\hat{h} = (\hat{t}, \hat{p}_M, \hat{M})$, the history on which we are changing the recommendations. We use h to denote a typical element of $T \times \mathcal{M} \times \mathcal{M}$.

Define the new mechanism, $(\gamma_{\mathcal{M}}^*, \gamma_{\mathcal{L}}^*, \gamma_X^*)$, as follows. First, $\gamma_{\mathcal{M}}^* = \gamma_{\mathcal{M}}$. Let $\gamma_{\mathcal{L}}^*$ satisfy $\gamma_{\mathcal{L}}^*(h)(m) = \gamma_{\mathcal{L}}(h)(m)$ if $h \neq \hat{h}$. Similarly, let $\gamma_{\mathcal{L}}^*(\hat{h})(m) = \gamma_{\mathcal{L}}(\hat{h})(m)$ for $m \notin \{\hat{m}, m_M^*\}$. Finally, let

$$\gamma_{\mathcal{L}}^*(\hat{h})(m) = \begin{cases} \gamma_{\mathcal{L}}(\hat{h})(m_M^*) + \gamma_{\mathcal{L}}(\hat{h})(\hat{m}), & \text{if } m = m_M^*; \\ 0, & \text{if } m = \hat{m}. \end{cases}$$

In other words, the probability that was on recommendation \hat{m} is moved to m_M^* .

Let $\gamma_X^*(h, m, m')(x) = \gamma_X(h, m, m')(x)$ if $(h, m) \neq (\hat{h}, m_M^*)$. In other words, on histories other than \hat{h} and on \hat{h} if the principal did not request maximal evidence, we do not change the mechanism's outcome. Also, for all $m \in \mathcal{L} \setminus \{m_M^*\}$, we set $\gamma_X^*(\hat{h}, m_M^*, m)(x)$ equal to

$$\frac{\gamma_{\mathcal{L}}(\hat{h})(\hat{m})\gamma_X(\hat{h}, \hat{m}, m)(x) + \gamma_{\mathcal{L}}(\hat{h})(m_M^*)\gamma_X(\hat{h}, m_M^*, m)(x)}{\gamma_{\mathcal{L}}(\hat{h})(\hat{m}) + \gamma_{\mathcal{L}}(\hat{h})(m_M^*)}.$$

Finally, we set $\gamma_X^*(\hat{h}, m_M^*, m_M^*)(x)$ equal to

$$\frac{\gamma_{\mathcal{L}}(\hat{h})(\hat{m})\gamma_X(\hat{h}, \hat{m}, \hat{m})(x) + \gamma_{\mathcal{L}}(\hat{h})(m_M^*)\gamma_X(\hat{h}, m_M^*, m_M^*)(x)}{\gamma_{\mathcal{L}}(\hat{h})(\hat{m}) + \gamma_{\mathcal{L}}(\hat{h})(m_M^*)}.$$

In other words, if m_M^* is requested and anything else is reported, then the response is the ‘‘average response’’ to this form of disobedience, averaging over the cases where \hat{m}

or $m_{\hat{M}}^*$ was requested in the original mechanism. On the other hand, if $m_{\hat{M}}^*$ is requested and reported, then the response is the average response to obedience in response to a request for either \hat{m} or $m_{\hat{M}}^*$ in the original mechanism.

We first show that this change in the mechanism does not change the outcome if the agent is truthful and obedient. The only situation a truthful and obedient agent is affected by the change is when her type is \hat{t} , the principal recommends (and she chooses) action \hat{a} , and the resulting message set is \hat{M} . Conditional on history \hat{h} and obeying the principal's recommendations, the probability of x in the new mechanism is

$$\begin{aligned}
& \sum_{m \in \mathcal{L}} \gamma_{\mathcal{L}}^*(\hat{h})(m) \gamma_X^*(\hat{h}, m, m)(x) \\
&= \sum_{m \in \mathcal{L} \setminus \{\hat{m}, m_{\hat{M}}^*\}} \gamma_{\mathcal{L}}(\hat{h})(m) \gamma_X(\hat{h}, m, m)(x) \\
&\quad + 0 + \gamma_{\mathcal{L}}^*(\hat{h})(m_{\hat{M}}^*) \gamma_X^*(\hat{h}, m_{\hat{M}}^*, m_{\hat{M}}^*)(x) \\
&= \sum_{m \in \mathcal{L} \setminus \{\hat{m}, m_{\hat{M}}^*\}} \gamma_{\mathcal{L}}(\hat{h})(m) \gamma_X(\hat{h}, m, m)(x) \\
&\quad + [\gamma_{\mathcal{L}}(\hat{h})(\hat{m}) + \gamma_{\mathcal{L}}(\hat{h})(m_{\hat{M}}^*)] \gamma_X^*(\hat{h}, m_{\hat{M}}^*, m_{\hat{M}}^*)(x) \\
&= \sum_{m \in \mathcal{L}} \gamma_{\mathcal{L}}(\hat{h})(m) \gamma_X(\hat{h}, m, m)(x).
\end{aligned}$$

Hence, as asserted, the outcome under truth-telling is the same in the new mechanism as in the original mechanism. Therefore, the agent's expected payoff from truth-telling and obedience is the same in the two mechanisms.

We now show that for any type t and any deviation feasible for t in the new mechanism, there is a deviation that is feasible for type t in the original mechanism which yields the same expected payoff. Since truth-telling is superior to any feasible deviation in the original mechanism, then, truth-telling is superior to any feasible deviation in the new mechanism.

To see this, fix any type t (which may equal \hat{t}) and consider any feasible deviation. Obviously, if the deviation involves reporting a type other than \hat{t} , this deviation is also available in the original mechanism and yields the same payoff in the new mechanism as in the original one since the way the mechanism responds to such a report has not changed. Hence we can restrict attention to deviations which involve reporting type \hat{t} . So fix any such deviation. Clearly, we may as well condition on the event that the principal requests the distribution \hat{p}_M , the agent chooses p_M (which may equal \hat{p}_M), the agent obtains message set M , and reports message set \hat{M} (which may equal M). Let $z : \hat{M} \rightarrow M$ give the message the agent sends as a function of the message the principal

requests from her. Then the agent's expected payoff conditional on this event is

$$\sum_{(x,m) \in X \times \mathcal{L}} \gamma_{\mathcal{L}}^*(\hat{h})(m) \gamma_X^*(\hat{h}, m, z(m))(x) u(t, p_M, x).$$

We can write this as

$$\begin{aligned} & \sum_{(x,m) \in X \times (\mathcal{L} \setminus \{\hat{m}, m_M^*\})} \gamma_{\mathcal{L}}(\hat{h})(m) \gamma_X(\hat{h}, m, z(m))(x) u(t, p_M, x) \\ & + \gamma_{\mathcal{L}}^*(\hat{h})(m_M^*) \sum_{x \in X} \gamma_X^*(\hat{h}, m_M^*, z(m_M^*))(x) u(t, p_M, x). \end{aligned}$$

We have two cases. First, suppose $z(m_M^*) \neq m_M^*$. In this case, the last term is equal to

$$\sum_{(x,m) \in X \times \{\hat{m}, m_M^*\}} \gamma_M(\hat{h})(m) \gamma_X(\hat{h}, m, z(m_M^*))(x) u(t, p_M, x).$$

Thus the conditional payoff to the deviation in the new mechanism is the same as the conditional payoff in the original mechanism where the agent responds to a request for *either* \hat{m} or m_M^* by sending $z(m_M^*)$. So in this case, the payoff to the deviation in the new mechanism is the same as the payoff to a certain deviation which was also feasible in the original mechanism.

Second, suppose $z(m_M^*) = m_M^*$. In this case, the last term is equal to

$$\sum_{(x,m) \in X \times \{\hat{m}, m_M^*\}} \gamma_M(\hat{h})(m) \gamma_X(\hat{h}, m, m)(x) u(t, p_M, x).$$

In other words, the payoff in the new mechanism is the same as the payoff in the old mechanism where the agent responds to a request for \hat{m} with \hat{m} and a request for m_M^* with m_M^* . Note that we are assuming that the deviation in the new mechanism is feasible for the agent, so $m_M^* \in M$. By the definition of normality, this implies $\hat{m} \in M$. Hence this deviation has the same payoff as a feasible deviation in the original mechanism.

In either case, then, the best deviation payoff in the new mechanism cannot exceed the best deviation payoff in the original mechanism, so the new mechanism is incentive compatible.

Clearly, we can repeat this argument as needed to obtain an incentive compatible mechanism which has the same mechanism outcome as γ and which has the property that $\gamma_{\mathcal{L}}(t, p_M, M)(m_M^*) = 1$ for all $(t, p_M, M) \in T \times \mathcal{M} \times \mathcal{M}$.

B Details for Example 1

Recall that there are two types, each with prior probability $1/2$, and each type has only one possible probability distribution over evidence sets. Type t_1 obtains evidence set $\{m_1\}$ with probability $1/3$, $\{m_2\}$ with probability $1/3$, and $\{m_1, m_2\}$ with probability $1/3$. Type t_2 receives evidence set $\{m_2\}$ with probability 1.

We give two examples of preferences, one for which the principal is strictly better off requesting message m_1 when the agent reports evidence set $\{m_1, m_2\}$ and one where he is strictly better off requesting message m_2 .

Before giving the examples, we note that we can simplify mechanisms for this setting in some innocuous ways. First, there is no need for the principal to recommend a choice of a distribution over evidence since each type has only one such distribution available. Similarly, if the agent reports a singleton evidence set, there is no need for the principal to recommend a particular message. Thus a mechanism starts with a type report by the agent. If the agent reports t_1 , she then makes a report of the realized evidence set. If this set is $\{m_1, m_2\}$, the principal then recommends a message to send from this set. Finally, after the agent sends a message, there is an outcome as a function of all that the principal has seen. Incentive compatibility means that the agent reports truthfully and follows the principal's recommendation if she reports evidence set $\{m_1, m_2\}$.

First, we present preferences for which the principal is strictly better off requesting m_2 from $\{m_1, m_2\}$. For this example, let $X = \{x, y, z\}$. Because there is only one possible distribution over evidence sets for each type, we can write utility for the principal and the agent as a function of x and t only. Assume the agent's utility as a function of her type is given by

	t_1	t_2
x	1	0
y	0	1
z	2	1

while the principal's utility is given by

	t_1	t_2
x	0	M
y	1	0
z	-2	0

where $M > 0$ and large. We write a typical lottery over X in the form (q_x, q_y, q_z) . We also write x interchangeably with $(1, 0, 0)$ and similarly for other degenerate lotteries.

Consider the following mechanism. Regardless of the type/evidence set reports by the agent and the requested message (if any) by the principal, if the message sent by the

agent is m_2 , the outcome is x . Regardless of the reports and requested message, if the message sent by the agent is m_1 , the outcome is y . The principal requests message m_2 when the agent reports type t_1 and message set $\{m_1, m_2\}$.

It is not hard to see that the agent may as well report truthfully and will obey the principal's recommendation. If the agent is type t_2 , the only evidence message she can send is m_2 , which necessarily generates outcome x . Hence she may as well report her type truthfully. If the agent is type t_1 and has evidence set $\{m_1\}$ or evidence set $\{m_2\}$, again her report of an evidence set does not affect the outcome, so she may as well report the evidence set truthfully. If her evidence set is $\{m_1, m_2\}$, she will want to send message m_2 so she will follow the principal's recommendation and has no incentive to misreport the set of evidence she has. Finally, t_1 has no incentive to misreport her type at the outset since the type report itself does not affect the outcome she will generate.

It is easy to see that this mechanism yields the principal an expected payoff of $(1/2)M + (1/6)$.

Now suppose the principal requests m_1 when the agent reports type t_1 and evidence set $\{m_1, m_2\}$. Without loss of generality, we can assume that if the principal observes a history inconsistent with obedience and truth-telling where the agent sends evidence message m_1 , then he chooses y . This is because m_1 proves that the type is t_1 and y is the worst outcome for t_1 . Hence this does the best possible job of preventing such deviations. We also may as well assume that there is some lottery δ which is the principal's response to any history inconsistent with obedience and truth-telling where the agent sends evidence message m_2 . This is because all such deviations are feasible for any type and hence all are required to be worse for any type than truth-telling and obedience. Since the outcomes from these various situations have the same direct payoff consequences for the principal (none since they are off path) and are subject to the same constraints, there is no need to make them different.

Let α^i , $i = 1, 2$, denote the outcome when the agent reports type t_1 , evidence set $\{m_i\}$, and sends evidence message m_i . Let α^3 denote the outcome when the agent reports t_1 , evidence set $\{m_1, m_2\}$, and obeys the principal recommendation to send m_1 . Finally, let β denote the outcome when the agent reports t_2 and sends evidence message m_2 .

To write the incentive compatibility constraints succinctly, let \succeq_i denote the preference relation for type t_i , $i = 1, 2$. We require that each type reports truthfully and obeys rather than doing some obvious deviation involving m_2 (e.g., reporting a type and/or message set that implies she must have a certain evidence message and then sending the other). Hence incentive compatibility requires $\alpha^2 \succeq_1 \delta$ (otherwise, when t_1 has $\{m_2\}$, she could claim evidence set $\{m_1\}$), $\alpha^3 \succeq_1 \delta$ (so when t_1 has evidence set $\{m_1, m_2\}$, she has no incentive to claim $\{m_1\}$ and the prove she lied), and $\beta \succeq_2 \delta$ (so t_2 does not lie and claim to be t_1 with evidence set $\{m_1\}$). Note that t_1 never has an incentive to lie in an

obvious way and send m_1 since this leads to her least preferred outcome y .

There are only a few situations in which the agent can misreport without being automatically caught. First, if she is type t_2 , she could claim to be t_1 and to have message set $\{m_2\}$. Hence we require $\beta \succeq_2 \alpha^2$. Second, if the agent is type t_1 and has evidence set $\{m_1, m_2\}$, she could claim either evidence set $\{m_1\}$ or $\{m_2\}$ and provide the evidence message consistent with this. Hence we require $\alpha^3 \succeq_1 \alpha^1, \alpha^2$. Third, because the agent knows she will be asked to report m_1 if she claims to have $\{m_1, m_2\}$, an agent of type t_1 with evidence set $\{m_1\}$ can claim to have evidence set $\{m_1, m_2\}$. Hence we require $\alpha^1 \succeq_1 \alpha^3$.

The only other constraint is that the agent reports her type honestly at the outset. For t_2 , this is implied by $\beta \succeq_2 \delta$ and $\beta \succeq_2 \alpha^2$. For t_1 , we require that $(1/3)(\alpha^1 + \alpha^2 + \alpha^3) \succeq_1 (1/3)y + (2/3)\beta$.

Summarizing, the constraints are

$$\begin{aligned} \alpha^1 &\sim_1 \alpha^3 \succeq_1 \alpha^2 \succeq_1 \delta \\ \frac{1}{3}\alpha^1 + \frac{1}{3}\alpha^2 + \frac{1}{3}\alpha^3 &\succeq_1 \frac{1}{3}y + \frac{2}{3}\beta \\ \beta &\succeq_2 \alpha^2, \delta. \end{aligned}$$

From here, we see that we may as well set $\delta = \alpha^2$ since this will satisfy the constraints and, since δ is off path, has no direct payoff consequences for the principal. Also, there is no need to distinguish between α^1 and α^3 since they enter entirely symmetrically. So we may as well take $\alpha^3 = \alpha^1$. This reduces the constraints to

$$\begin{aligned} \alpha^1 &\succeq_1 \alpha^2 \\ \frac{2}{3}\alpha^1 + \frac{1}{3}\alpha^2 &\succeq_1 \frac{1}{3}y + \frac{2}{3}\beta \\ \beta &\succeq_2 \alpha^2. \end{aligned}$$

Rewriting in terms of expected utilities and doing some rearranging gives

$$\alpha_x^1 + 2\alpha_z^1 \geq \alpha_x^2 + 2\alpha_z^2 \tag{1}$$

$$\alpha_x^1 + 2\alpha_z^1 + \frac{1}{2}[\alpha_x^2 + 2\alpha_z^2] \geq \beta_x + 2\beta_z \tag{2}$$

$$\beta_x \leq \alpha_x^2. \tag{3}$$

It is not hard to see that we may as well set $\beta_z = 0$. If it is strictly positive, we can reduce it to zero, increasing β_y to compensate. Since the principal gets 0 utility from y

or z given type t_2 , this has no direct payoff consequences but may relax the constraints. Hence we can rewrite equation (2) as

$$\alpha_x^1 + 2\alpha_z^1 + \frac{1}{2}[\alpha_x^2 + 2\alpha_z^2] \geq \beta_x.$$

But we have

$$\alpha_x^1 + 2\alpha_z^1 + \frac{1}{2}[\alpha_x^2 + 2\alpha_z^2] \geq \alpha_x^1 + 2\alpha_z^1 \geq \alpha_x^2 + 2\alpha_z^2 \geq \alpha_x^2 \geq \beta_x$$

where the first inequality is from non-negativity of the probabilities, the second from equation (1), the third from non-negativity, and the fourth from equation (3). Hence constraint (2) is redundant.

From the remaining constraints, then, we see that we must have $\beta_x = \alpha_x^2$. This is because the principal wants β_x as large as possible and it affects no constraints other than (3). Finally, then, we can write the principal's problem as maximizing

$$\frac{1}{2}\alpha_x^2 M + \frac{1}{3}[\alpha_y^1 - 2\alpha_z^1] + \frac{1}{6}[\alpha_y^2 - 2\alpha_z^2]$$

subject to equation (1). If we substitute for α_y^i using $\alpha_x^i + \alpha_y^i + \alpha_z^i = 1$ and rearrange, we can write the objective function as a constant plus

$$-\frac{1}{3}\alpha_x^1 - \alpha_z^1 - \frac{1}{2}\alpha_z^2 + \alpha_x^2 \left(\frac{1}{2}M - \frac{1}{6} \right).$$

Clearly, the remaining constraint (1) must bind — otherwise, $\alpha_x^1 = \alpha_z^1 = \alpha_z^2 = 0$ and $\alpha_x^2 = 1$, contradicting the constraint. Hence

$$\alpha_x^1 = \alpha_x^2 + 2\alpha_z^2 - 2\alpha_z^1.$$

Substituting into the objective function gives

$$-\frac{1}{3}\alpha_z^1 - \frac{7}{6}\alpha_z^2 + \alpha_x^2 \left(\frac{1}{2}M - \frac{1}{2} \right).$$

Assuming $M > 1$, the solution is $\alpha_x^2 = 1$ and $\alpha_z^2 = \alpha_z^1 = 0$, implying $\alpha_x^1 = 1$ as well.

In short, the optimal mechanism which requests m_1 when the agent reports evidence set $\{m_1, m_2\}$ has outcome x regardless of the agent's type or evidence set. This gives the principal a payoff of $M/2$, strictly below what he obtained from the mechanism which requested m_2 from $\{m_1, m_2\}$.

The preferences for which the principal is strictly better off requesting m_1 are simpler. Here we let $X = \{x, y\}$, take the agent's utility function to be

	t_1	t_2
x	0	0
y	1	1

and assume the principal's utility is given by

	t_1	t_2
x	0	M
y	1	0

where $M > 0$ and large. First, consider the following mechanism where the principal induces the agent to send m_1 from evidence set $\{m_1, m_2\}$: the principal chooses y whenever the message sent is m_1 and x whenever it is m_2 . Since both types of agents prefer y , the agent will send m_1 whenever possible. So the principal will get the outcome x when the type is t_2 and get y with probability $1/3$, x otherwise, when the type is t_1 . The principal's payoff is $(1/2)M + (1/3)$.

Next consider the best mechanism for the principal with the property that he induces the agent to send m_2 from $\{m_1, m_2\}$. As before, let α^i be the lottery chosen by the principal when the agent reports t_1 , reports evidence set $\{m_i\}$, and sends message m_i . Let α^3 be the outcome when the agent reports t_1 and evidence set $\{m_1, m_2\}$ and then follows the principal's recommendation to set message m_2 . Let β be the outcome when the type report is t_2 , the evidence set report is $\{m_2\}$, and the message is m_2 . Let δ^i denote the outcome when the history is clearly off path and the message sent is m_i . Then the incentive compatibility constraints are as follows. Type t_1 with evidence set $\{m_i\}$ must prefer α^i to the off path deviation she can generate of δ^i , $i = 1, 2$. Similarly, $\beta \succeq_2 \delta^2$. For undetectable deviations, we require $\alpha^2 \succeq_1 \alpha^3$ since t_1 with evidence set $\{m_2\}$ can pretend to have evidence set $\{m_1, m_2\}$. Also, we require $\alpha^3 \succeq_1 \alpha^1, \alpha^2$. Since t_2 could pretend to be t_1 with either $\{m_2\}$ or $\{m_1, m_2\}$, we require $\beta \succeq_2 \alpha^2, \alpha^3$. Finally, we require

$$\frac{1}{3}\alpha^1 + \frac{1}{3}\alpha^2 + \frac{1}{3}\alpha^3 \succeq_1 \frac{1}{3}\delta^1 + \frac{2}{3}\beta.$$

It is not hard to see that we may as well set $\alpha^2 = \delta^2$ since δ^2 has no direct payoff consequences and this will necessarily satisfy the constraints. Also, given the symmetry of the way they enter the problem, we may as well set $\alpha^2 = \alpha^3$. Finally, given that the constraints are weaker as δ^1 is made worse for t_1 , we may as well set $\delta^1 = x$.

Using this to simplify, the constraints are

$$\begin{aligned} \alpha^2 &\succeq_1 \alpha^1 \\ \frac{1}{3}\alpha^1 + \frac{2}{3}\alpha^2 &\succeq_1 \frac{1}{3}x + \frac{2}{3}\beta \\ \beta &\succeq_2 \alpha^2. \end{aligned}$$

Rewriting using the agent's utility function, we have

$$\alpha_y^2 \geq \alpha_y^1$$

$$\frac{1}{3}\alpha_y^1 + \frac{2}{3}\alpha_y^2 \geq \frac{2}{3}\beta_y$$

$$\beta_y \geq \alpha_y^2.$$

Since the principal prefers x to y when the agent is type t_2 , the last constraint must bind, so $\beta_y = \alpha_y^2$. This implies that the middle constraint does not bind. Hence we can reduce the problem to maximizing

$$\frac{1}{2}(1 - \alpha_y^2)M + \frac{1}{2}\left[\frac{1}{3}\alpha_y^1 + \frac{2}{3}\alpha_y^2\right]$$

subject to $\alpha_y^2 \geq \alpha_y^1$. It is easy to see that if M is large enough, the solution is $\alpha_y^1 = \alpha_y^2 = 0$.

Hence for M large, the best mechanism for the principal which induces the agent to send m_2 from evidence set $\{m_1, m_2\}$ is the constant mechanism x , which gives the principal the payoff $M/2$. So the principal is strictly worse off with such a mechanism.

C Proof of Corollary 1

Fix an incentive compatible mechanism $\gamma = (\gamma_{\mathcal{M}}, \gamma_{\mathcal{L}}, \gamma_X)$. By Theorem 1, we can assume without loss of generality that $\gamma_{\mathcal{L}}(t, p_M, M)(m_M^*) = 1$ for all $(t, p_M, M) \in T \times \mathcal{M} \times \mathcal{M}$. We construct a mechanism $(\gamma_{\mathcal{M}}^*, \gamma_X^*)$ for the abbreviated protocol which is incentive compatible and has the same outcome as γ . To do so, first let $\gamma_{\mathcal{M}}^* = \gamma_{\mathcal{M}}$.

To construct γ_X^* , note that in the abbreviated protocol, $\gamma_X^* : T \times \mathcal{M} \times \mathcal{L} \rightarrow \Delta(X)$, while in the full protocol, $\gamma_X : T \times \mathcal{M} \times \mathcal{M} \times \mathcal{L} \times \mathcal{L} \rightarrow \Delta(X)$ since the choice of x can depend on the agent's report of an evidence set and the message the principal requests, in addition to the type report, distribution recommendation, and received message as in the abbreviated protocol.

Given any $m \in \mathcal{L}$, define $M^*(m)$ as follows. First, if there is any M such that $m = m_M^*$, then let $M^*(m)$ equal this message set M .¹⁴ Otherwise, let $M^*(m)$ denote any $M \in \mathcal{M}$ such that $m \in M$. Given this, let

$$\gamma_X^*(t, p_M, m) = \gamma_X(t, p_M, M^*(m), m_{M^*(m)}^*, m).$$

In other words, if the agent reports t , the principal recommends p_M , and the agent shows message m , then the outcome is the same as in the original mechanism when the agent reports t , the principal recommends p_M , the agent reports evidence set $M^*(m)$,

¹⁴It is straightforward to show that if $m_M^* = m_{\hat{M}}^*$, then $M = \hat{M}$. That is, $M^*(m)$ is unambiguously defined in this case.

the principal requests the maximal evidence message for this set, and the agent provides message m .

If the agent truthfully reports her type, follows the principal's recommended distribution p_M , and provides the maximal evidence message from any evidence set she obtains, this construction implies that the resulting distribution over X in the new mechanism will be the same as in the original mechanism. Hence if this mechanism is incentive compatible, it yields the same outcome as the original mechanism.

So consider an agent of type t who reports type \hat{t} (which may or may not equal t), has p_M recommended to her by the principal, chooses \hat{p}_M , obtains evidence set M , and sends message m from it. In this situation, the outcome under the new mechanism is $\gamma_X(\hat{t}, p_M, M^*(m), m_{M^*(m)}^*, m)$, exactly the same outcome the agent could have obtained by reporting \hat{t} , choosing \hat{p}_M , reporting $M^*(m)$ as her evidence set, and then sending m . That is, any outcome the agent can generate in the new mechanism using a strategy which deviates from truth-telling, obedience, and sending maximal evidence is an outcome she could have generated in the original mechanism using a certain strategy which deviated from truth-telling and obedience. Since the original mechanism was incentive compatible, truth-telling and obedience were superior to this deviation. Hence the agent prefers truth-telling, obedience, and maximal evidence in the new mechanism to any deviation, so the mechanism is incentive compatible.

D Proof of Theorem 2

Fix an incentive compatible mechanism for the evidence-acquisition model under normality. By Corollary 1, we can take this mechanism to be based on the abbreviated protocol. Hence it consists of a pair of functions $\gamma_{\mathcal{M}} : T \rightarrow \Delta(\mathcal{M})$ and $\gamma_X : T \times \mathcal{M} \times \mathcal{L} \rightarrow \Delta(X)$. For the signal choice model, a mechanism is a pair of functions $\gamma_S^* : T \rightarrow \Delta(S)$ and $\gamma_X^* : T \times S \times \mathcal{L} \rightarrow \Delta(X)$.

Given the incentive compatible mechanism for the abbreviated protocol, we construct an equivalent incentive compatible mechanism for the associated signal-choice model as follows. Let

$$\gamma_S^*(t)(s_{(p_M, \sigma^*)}) = \gamma_{\mathcal{M}}(t)(p_M).$$

That is, given a report of t , the principal recommends the signal distribution generated by evidence distribution p_M followed by showing maximal evidence with the same probability he recommended p_M in the original mechanism. Let

$$\gamma_X^*(t, s_{(p_M, \sigma^*)}, m) = \gamma_X(t, p_M, m).$$

That is, if the agent report type t and the signal distribution the principal recommends is the one corresponding to p_M and maximal evidence, then the principal replies to message m in the new mechanism the same way he replied in the original mechanism given type report t and recommendation p_M .

It is easy to see that if the agent reports her type truthfully and follows the principal's recommended signal distribution, then the outcome is equivalent to that of the original mechanism as defined in the statement of the theorem. If the agent deviates, this corresponds directly to a particular deviation strategy in the original mechanism and hence cannot be profitable for her. In particular, if type t reports \hat{t} , receives the recommendation s_{p_M, σ^*} , and uses signal distribution $s_{(\hat{p}_M, \hat{\sigma})}$ instead, she generates exactly the outcome she would have generated in the original mechanism if she reported \hat{t} , received the recommendation p_M , chose the distribution \hat{p}_M instead, and selected a message to send using the function $\hat{\sigma}$. Hence the mechanism is incentive compatible.

E Proof of Theorem 3

Fix an incentive compatible mechanism (γ_S, γ_X) where $\gamma_S(t_1)(\hat{s}_1) = \hat{\alpha} > 0$. Let $\alpha = \gamma_S(t_1)(s_1)$ (where this can be 0). We construct an incentive compatible mechanism (γ_S^*, γ_X^*) with the same outcome where the principal recommends s_1 to t_1 with probability $\alpha + \hat{\alpha}$ and never recommends \hat{s}_1 to t_1 .

For any $t \neq t_1$, $\gamma_S^*(t) = \gamma_S(t)$ and $\gamma_X^*(t, s, m) = \gamma_X(t, s, m)$ for all (s, m) . For $s \neq s_1, \hat{s}_1$, we have $\gamma_S^*(t_1)(s) = \gamma_S(t_1)(s)$ and $\gamma_X^*(t_1, s, m) = \gamma_X(t_1, s, m)$. That is, if the agent reports a type other than t_1 , the new mechanism is the same as the original one and if the agent reports t_1 , the principal recommends signals other than s_1 or \hat{s}_1 with the same probability and treats them the same way as in the original mechanism.

Let $\gamma_S^*(t_1)(\hat{s}_1) = 0$ and $\hat{\gamma}_S^*(t_1)(s_1) = \alpha + \hat{\alpha}$. Since the principal never recommends \hat{s}_1 in response to a report of t_1 in this mechanism, we only need to specify $\gamma_X^*(t, s, m)$ for $(t, s) = (t_1, s_1)$. For notational convenience, we enumerate the messages as $\mathcal{L} = \{m_1, \dots, m_L\}$ and for the Markov matrix Λ , we write the entry corresponding to (m_i, m_j) as λ_{ij} rather than λ_{m_i, m_j} .

Let

$$\gamma_X^*(t_1, s_1, m_i) = \frac{\alpha}{\alpha + \hat{\alpha}} \gamma_X(t_1, s_1, m_i) + \frac{\hat{\alpha}}{\alpha + \hat{\alpha}} \sum_j \lambda_{ij} \gamma_X(t_1, \hat{s}_1, m_j).$$

Because all the λ_{ij} 's are non-negative and because $\sum_j \lambda_{ij} = 1$ for every i , we see that $\gamma_X^*(t_1, s_1, m_i)$ is a convex combination of probability distributions over X and hence is a

probability distribution over X .

Given this specification, suppose all types report honestly and obey the principal's recommendations. Obviously, if the true type $t \neq t_1$, we have the same outcome as before. So suppose $t = t_1$. Then the expected outcome is

$$(\alpha + \hat{\alpha}) \sum_i s_1(m_i) \gamma_X^*(t_1, s_1, m_i) + \sum_{s \in S_{t_1} \setminus \{s_1, \hat{s}_1\}} \gamma_S^*(t_1)(s) \sum_M s(m) \gamma_X^*(t_1, s, m). \quad (4)$$

Substituting for γ_X^* , the first term in equation (4) is

$$\begin{aligned} & \alpha \sum_i s_1(m_i) \gamma_X(t_1, s_1, m_i) + \hat{\alpha} \sum_i s_1(m_i) \sum_j \lambda_{ij} \gamma_X(t_1, \hat{s}_1, m_j) \\ & = \alpha \sum_i s_1(i) \gamma_X(t_1, s_1, m_i) + \hat{\alpha} \sum_j \gamma_X(t_1, \hat{s}_1, m_j) \sum_i s_1(m_i) \lambda_{ij}. \end{aligned}$$

But $s_1 \Lambda = \hat{s}_1$, so that for every j , $\sum_i s_1(m_i) \lambda_{ij} = \hat{s}_1(m_j)$. Hence this is

$$= \alpha \sum_i s_1(m_i) \gamma_X(t_1, s_1, m_i) + \hat{\alpha} \sum_i \hat{s}_1(m_i) \gamma_X(t_1, \hat{s}_1, m_j).$$

Substituting this for the first term in equation (4) and substituting for γ_S^* and γ_X^* in the second term, we see that the expected outcome under truth-telling and obedience is the same as under the original mechanism.

To show that the new mechanism is incentive compatible, we show that any deviation from truth-telling and obedience by any type generates a distribution over outcomes that the same type could have generated in the original mechanism. Since the original mechanism was incentive compatible, this deviation is not profitable, so the new mechanism is incentive compatible.

To see that this holds, fix any type t and any signal $s' \in S_t$. If t makes any type report other than t_1 , the mechanism has not changed, so the claim obviously holds. So suppose type t reports type t_1 . If the mechanism makes any signal recommendation other than s_1 , then, again, the mechanism is the same as before, so the claim holds. So suppose the mechanism recommends signal s_1 and the agent uses s' . The expected outcome times the probability of this event is

$$(\alpha + \hat{\alpha}) \sum_i s'(m_i) \gamma_X^*(t_1, s_1, m_i) = \alpha \sum_i s'(m_i) \gamma_X(t_1, s_1, m_i) + \hat{\alpha} \sum_i s'(m_i) \sum_j \lambda_{ij} \gamma_X(t_1, \hat{s}_1, m_j).$$

By assumption, $s' \Lambda \in \text{conv}(S_t)$. Hence we can write $s' \Lambda = \sum_k a_k s^k$ where $a_k \geq 0$ for all k , $\sum_k a_k = 1$, and $s^k \in S_t$ for all k . In particular, for every j ,

$$\sum_i s'(m_i) \lambda_{ij} = \sum_k a_k s^k(m_j).$$

Hence we can rewrite the above as

$$\alpha \sum_i s'(i) \gamma_X(t_1, s_1, m_i) + \hat{\alpha} \sum_k a_k s^k(i) \gamma_X(t_1, \hat{s}_1, m_i).$$

This is exactly what t would generate in the original mechanism if she responded to a recommendation of s_1 with s' and a recommendation of \hat{s}_1 by randomizing with probability a_k on s^k . Thus, as asserted, any expected outcome t can generate in the new mechanism is identical to some outcome she could have generated in the original mechanism. Hence the new mechanism is incentive compatible.

F Proof of Theorem 5

F.1 Lemma

The following result will be useful. Let W be a finite set of states of the world and A a finite set of actions. Let $u : A \times W \rightarrow \mathbf{R}$ be a utility function. Say that $\sigma \in \Delta(A)$ is a *best reply to* $p \in \Delta(W)$ if

$$\sum_w p(w) \sum_a \sigma(a) u(a, w) \geq \sum_w p(w) \sum_a \sigma'(a) u(a, w), \quad \forall \sigma' \in \Delta(A).$$

Say that σ is a *best reply* if there exists $p \in \Delta(W)$ such that σ is a best reply to p . Say that σ is *strictly dominated* if there exists $\sigma' \in \Delta(A)$ such that

$$\sum_a \sigma'(a) u(a, w) > \sum_a \sigma(a) u(a, w), \quad \forall w \in W.$$

Standard results say that σ is a best reply if and only if it is not strictly dominated. (This is typically stated for pure strategies σ , but it applies to mixed as well.)

Lemma 1. *Suppose σ' is strictly dominated. Then there exists a mixed strategy $\hat{\sigma}$ which strictly dominates σ' and which is not itself strictly dominated.*

Proof. Suppose not. That is, suppose σ' is strictly dominated, but is not strictly dominated by any undominated mixed strategy. Let Σ^* denote the set of undominated strategies in $\Delta(A)$. Equivalently, Σ^* is the set of all $\sigma \in \Delta(A)$ that are best replies. By finiteness of A , this set is nonempty.

Let

$$\mathcal{U} = \text{conv} \left(\left\{ u \in \mathbf{R}^W \mid \exists \sigma \in \Sigma^* \cup \{\sigma'\} \text{ with } u_w = \sum_a \sigma(a) u(a, w), \quad \forall w \right\} \right).$$

$$\mathcal{U}^D = \left\{ u \in \mathbf{R}^W \mid u_w \geq \sum_a \sigma'(a)u(a, w), \forall w \right\}.$$

By hypothesis, there is no mixed strategy in Σ^* which strictly dominates σ' . Hence $\mathcal{U} \cap \text{int}(\mathcal{U}^D) = \emptyset$, so the interiors of \mathcal{U} and \mathcal{U}^D are disjoint. Clearly, both sets are nonempty and convex. Hence there exists a separating hyperplane. That is, there is $p \in \mathbf{R}^W$ such that $p \neq 0$ and $p \cdot u \geq p \cdot \hat{u}$ for all $u \in \mathcal{U}^D$, $\hat{u} \in \mathcal{U}$.

Consider \hat{u} defined by $\hat{u}_w = \sum_a \sigma'(a)u(a, w)$. Obviously, this is an element of \mathcal{U} . Consider u defined by $u_w = \hat{u}_w$ for $w \neq w'$ and $u_{w'} = \hat{u}_{w'} + \varepsilon$ for some $\varepsilon > 0$ and some w' . Clearly, this is an element of \mathcal{U}^D . Hence the separating hyperplane satisfies $p_{w'}\varepsilon \geq 0$. Since w' is arbitrary, $p_w \geq 0$ for all w . Since $p \neq 0$, we can renormalize by replacing p with \hat{p} defined by $\hat{p}_w = p_w / \sum_{w'} p_{w'}$. Hence $\hat{p} \in \Delta(W)$.

Continuing with the same \hat{u} as above, we see that we have $\hat{p} \in \Delta(W)$ such that

$$\sum_w \hat{p}(w) \sum_a \sigma'(a)u(a, w) \geq \sum_w \hat{p}(w)u_w \quad \forall u \in \mathcal{U}.$$

In particular, for any $\sigma \in \Sigma^*$, we can let u be the vector defined by $u_w = \sum_a \sigma(a)u(a, w)$ to conclude that

$$\sum_w \hat{p}(w) \sum_a \sigma'(a)u(a, w) \geq \sum_w \hat{p}(w) \sum_a \sigma(a)u(a, w) \quad \forall \sigma \in \Sigma^*.$$

By hypothesis, σ' is strictly dominated by some mixed strategy, say $\hat{\sigma} \notin \Sigma^*$. Hence

$$\sum_w \hat{p}(w) \sum_a \hat{\sigma}(a)u(a, w) > \sum_w \hat{p}(w) \sum_a \sigma'(a)u(a, w) \geq \sum_w \hat{p}(w) \sum_a \sigma(a)u(a, w) \quad \forall \sigma \in \Sigma^*.$$

Hence no best reply to \hat{p} is contained in Σ^* , a contradiction. ■

F.2 Theorem 5

Fix the Nash equilibrium $(\hat{\beta}, \gamma^*)$ constructed in the proof of Theorem 4. We claim that without loss of generality, we can assume the agent's strategy is sequentially rational at all of her information sets. Obviously, it must be sequentially rational at every positive probability information set since $\hat{\beta}$ is a best reply to γ^* . So if it is not sequentially rational at some unreached information set, simply replace her strategy at any such unreached information set with any alternative strategy which is sequentially rational. This cannot lead the principal to deviate since

$$V(\hat{\beta}, \gamma^*) = \max_{\gamma \in \Gamma} \max_{\beta \in BR(\gamma)} V(\beta, \gamma).$$

So fix any unreached information set for the principal. First, suppose this information set is at a stage where the principal has observed only cheap talk statements by the agent, say r' . Let \hat{r} be any other profile of cheap talk statements which the principal could have heard with positive probability at this stage. Change the equilibrium so the principal's beliefs and continuation strategy after seeing r' are the same as those after seeing \hat{r} . Similarly, change the agent's continuation strategy after sending r' to match her continuation strategy after sending \hat{r} . It is easy to see that this must satisfy sequential rationality for the principal at this information set and does not attract a deviation from the equilibrium.

So, without loss of generality, we can assume that the unreached information set is one where the cheap talk the principal has observed from the agent, say r' , has positive probability but the joint observation of r' and the evidence message, say m' , does not. Since the protocol is very simple, this means that the principal chooses $x \in X$ at this information set. Let $\sigma^* \in \Delta(X)$ be the randomization chosen by the principal at this information set. Let T^* denote the set of t such that there exists $s \in S_t$ with $s(m') > 0$.

First, suppose there exists $\rho^* \in \Delta(T^*)$ such that σ^* maximizes over $\sigma \in \Delta(X)$

$$\sum_{t \in T^*} \rho^*(t) \sum_x \sigma(x) \nu(t) u(t, x).$$

In this case, take the principal's belief at this information set to be ρ^* . Then the principal's strategy is sequentially rational at this information set.

Second, suppose there is no $\rho^* \in \Delta(T^*)$ for which σ^* maximizes the principal's expected utility. By Lemma 1, σ^* is dominated with respect to T^* in the sense that there is some $\hat{\sigma} \in \Delta(X)$ such that

$$\sum_x \hat{\sigma}(x) \nu(t) u(t, x) > \sum_x \sigma^*(x) \nu(t) u(t, x), \quad \forall t \in T^*$$

and such that $\hat{\sigma}$ is not itself dominated in this sense. Since $\hat{\sigma}$ is not dominated in this sense, there exists $\hat{\rho} \in \Delta(T^*)$ such that $\hat{\sigma}$ maximizes the principal's expected utility. Set the principal's belief at this information set to equal $\hat{\rho}$ and his action at this information set to be $\hat{\sigma}$. So this satisfies sequential rationality for the principal. If this change in the principal's strategy does not make any deviation by any type of the agent profitable, we have a perfect Bayesian equilibrium.

So let \hat{T} denote the (nonempty) set of t for which this change in the principal's strategy does create a profitable deviation. It is easy to see that we must have $\hat{T} \subseteq T^*$ (since no other type could generate the information set in question) and for all $t \in \hat{T}$, $\sum_x \hat{\sigma}(x) u(t, x) > \sum_x \sigma^*(x) u(t, x)$. But recall that for all $t \in T^*$, we have we have $\nu(t) \sum_x \hat{\sigma}(x) u(t, x) > \nu(t) \sum_x \sigma^*(x) u(t, x)$, implying that $\nu(t) > 0$ for all $t \in \hat{T}$.

Let $\hat{\gamma}$ denote the strategy of the principal identical to γ^* *except* that the principal chooses $\hat{\sigma}$ on this information set. Let $\hat{\beta}'$ denote the best reply of the agent which differs from $\hat{\beta}$ only in letting types $t \in \hat{T}$ deviate. By hypothesis, $\hat{T} \neq \emptyset$, so $\hat{\beta}' \neq \hat{\beta}$.

Note that

$$\begin{aligned} V(\hat{\beta}', \hat{\gamma}) &= \mathbb{E}_t[\nu(t)U(\hat{\beta}', \hat{\gamma}, t)] \\ &= \Pr[\nu(t) < 0]\mathbb{E}_t[\nu(t)U(\hat{\beta}', \hat{\gamma}, t) \mid \nu(t) < 0] + \Pr[\nu(t) > 0]\mathbb{E}_t[\nu(t)U(\hat{\beta}', \hat{\gamma}, t) \mid \nu(t) > 0] \\ &> \Pr[\nu(t) < 0]\mathbb{E}_t[\nu(t)U(\hat{\beta}, \gamma^*, t) \mid \nu(t) < 0] + \Pr[\nu(t) > 0]\mathbb{E}_t[\nu(t)U(\hat{\beta}, \gamma^*, t) \mid \nu(t) > 0] \\ &= V(\hat{\beta}, \gamma^*) = V^*. \end{aligned}$$

The strict inequality comes from the fact that the types who deviate in response to $\hat{\gamma}$ are made strictly better off than they were at $(\hat{\beta}, \gamma^*)$. ■

G Multi-Agent Example

Assume $T_i = \{t_i^a, t_i^b\}$ where the types are equally likely. Assume $v_1(t_1^a) = 4$ and $v_1(t_1^b) = 1$, while $v_2(t_2^a) = 3$ and $v_2(t_2^b) = 0$. Assume each type of each agent has only one signal distribution available. Letting the unique signal choice available to t_i^k be denoted s_i^k , assume these distributions are

$$\begin{array}{cc|cc} & s_1^a & s_1^b & s_2^a & s_2^b \\ m_h & 1 & 3/4 & 1 & 0 \\ m_\ell & 0 & 1/4 & 0 & 1 \end{array}$$

So 2's type *must* be revealed to the principal. If 1 is type t_1^b , then her type is revealed with probability 1/4.

It is not hard to show that the optimal mechanism for our signal-choice protocol is:

$v_2(t_2)$	1's Report	1's Signal	Prob(1)
0	$v_1(t_1) = 1$	any	1
0	$v_1(t_1) = 4$	m_h	1
3	$v_1(t_1) = 1$	any	0
3	$v_1(t_1) = 4$	m_h	1/3
any	$v_1(t_1) = 4$	m_ℓ	0

In other words, when agent 2's type is revealed to be such that the value of giving the good to her is 0, then for any report and any signal realization consistent with this report from agent 1, the principal gives the good to agent 1. If agent 2's type is revealed to be such that the $v_2 = 3$ and agent 1 reports that $v_1 = 1$, the principal gives the good

to 2. If 2's type is revealed to be such that $v_2 = 3$, 1 claims that $v_1 = 4$, and the signal realization confirms this in the sense that the realization is m_h , the principal gives the good to 1 with probability $1/3$. Finally, whenever 1 is revealed to have lied, 2 gets the good.

There is no equilibrium with this outcome. The reason is that type reports cannot convey information in equilibrium. More specifically, note that if the signal realizations reveal the types, then sequential rationality for the principal dictates his allocation decision. Also, if 2's type is revealed to give 0 value to the principal, then sequential rationality dictates that he give the good to 1, regardless of any type report or signal realization for 1.

The only information sets where the principal's allocation is not uniquely determined by sequential rationality is when the value of giving the good to 2 is revealed to be 3 and the realization of 1's signal does not reveal her type. If the principal is more likely to give the good to 1 in this situation for one of the two possible type reports than for the other, then 1 will always use this type report. Hence the type report cannot convey any useful information.

It is not hard to use this to show that the only equilibrium outcome is that the principal gives the good to agent 2 if $v_2 = 3$ and gives it to 1 otherwise.

References

- [1] Acharya, V., P. DeMarzo, and I. Kremer, “Endogenous Information Flows and the Clustering of Announcements,” *American Economic Review*, **101**, December 2011, 2955–2979.
- [2] Ball, I., and D. Kattwinkel, “Probabilistic Verification in Mechanism Design,” working paper, August 2019.
- [3] Ben-Porath, E., E. Dekel, and B. Lipman, “Mechanisms with Evidence: Commitment and Robustness,” *Econometrica*, **87**, March 2019, 529–566.
- [4] Ben-Porath, E., and B. Lipman, “Implementation and Partial Provability,” *Journal of Economic Theory*, **147**, September 2012, 1689–1724.
- [5] Blackwell, D., and M. Girshick, *Theory of Games and Statistical Decisions*, Wiley, 1954.
- [6] Bull, J., and J. Watson, “Hard Evidence and Mechanism Design,” *Games and Economic Behavior*, **58**, January 2007, 75–93.
- [7] Chakraborty, A., and R. Harbaugh, “Persuasion by Cheap Talk,” *American Economic Review*, **100**, December 2010, 2361–2382.
- [8] Che, Y.-K., and N. Kartik, “Opinions as Incentives,” *Journal of Political Economy*, **117**, October 2009, 815–860.
- [9] Chen, Y., “Career Concerns and Excessive Risk Taking,” *Journal of Economics and Management Strategy*, **24**, Spring 2015, 110–130.
- [10] Crawford, V., and J. Sobel, “Strategic Information Transmission,” *Econometrica*, **50**, November 1982, 1431–1451.
- [11] Deb, R., M. Pai, and M. Said, “Evaluating Strategic Forecasters,” *American Economic Review*, **108**, October 2018, 3057–3103.
- [12] DeMarzo, P., I. Kremer, and A. Skrzypacz, “Test Design and Minimum Standards,” *American Economic Review*, **109**, June 2019, 2173–2207.
- [13] Deneckere, R. and S. Severinov, “Mechanism Design with Partial State Verifiability,” *Games and Economic Behavior*, **64**, November 2008, 487–513.
- [14] Dye, R. A., “Disclosure of Nonproprietary Information,” *Journal of Accounting Research*, **23**, 1985, 123–145.
- [15] Gerardi, D., and R. Myerson, “Sequential Equilibria in Bayesian Games with Communication,” *Games and Economic Behavior*, **60**, July 2007, 104–134.

- [16] Glazer, J., and A. Rubinstein, “On Optimal Rules of Persuasion,” *Econometrica*, **72**, November 2004, 1715–1736.
- [17] Glazer, J., and A. Rubinstein, “A Study in the Pragmatics of Persuasion: A Game Theoretical Approach,” *Theoretical Economics*, **1**, December 2006, 395–410.
- [18] Green, J., and J.-J. Laffont, “Partially Verifiable Information and Mechanism Design,” *Review of Economic Studies*, **53**, July 1986, 447–456.
- [19] Grossman, S. J., “The Informational Role of Warranties and Private Disclosures about Product Quality,” *Journal of Law and Economics*, **24**, 1981, 461–483.
- [20] Guttman, I., I. Kremer, and A. Skrzypacz, “Not Only What but also When: A Theory of Dynamic Voluntary Disclosure,” *American Economic Review*, **104**, August 2014, 2400–2420.
- [21] Halac, M., and I. Kremer, “Experimenting with Career Concerns,” Columbia Business School working paper, December 2017.
- [22] Hart, S., I. Kremer, and M. Perry, “Evidence Games: Truth and Commitment,” *American Economic Review*, **107**, March 2017, 690–713.
- [23] Hedlund, J., “Bayesian Persuasion by a Privately Informed Sender,” *Journal of Economic Theory*, **167**, January 2017, 229–268.
- [24] Holmstrom, B., “Managerial Incentive Problems: A Dynamic Perspective,” *Review of Economic Studies*, **66**, January 1999, 169–182.
- [25] Kamenica, E., and M. Gentzkow, “Bayesian Persuasion,” *American Economic Review*, **101**, October 2011, 2590–2615.
- [26] Kartik, N., and O. Tercieux, “Implementation with Evidence,” *Theoretical Economics*, **7**, May 2012, 323–355.
- [27] Kosenko, A., “Bayesian Persuasion with Private Information,” Columbia University working paper, February 2018.
- [28] Lipman, B., and D. Seppi, “Robust Inference in Communication Games with Partial Provability,” *Journal of Economic Theory*, **66**, August 1995, 370–405.
- [29] Lipnowski, E., and D. Ravid, “Cheap Talk with Transparent Motives,” working paper, May 2019.
- [30] Matthews, S., and L. Mirman, “Equilibrium Limit Pricing: The Effects of Private Information and Stochastic Demand,” *Econometrica*, **51**, July 1983, 981–996.
- [31] Matthews, S., and A. Postlewaite, “Quality Testing and Disclosure,” *RAND Journal of Economics*, **16**, Autumn 1985, 328–340.

- [32] Milgrom, P., “Good News and Bad News: Representation Theorems and Applications,” *Bell Journal of Economics*, **12**, 1981, 350–391.
- [33] Perez-Richet, E., “Interim Bayesian Persuasion: First Steps,” *American Economic Review: Papers & Proceedings*, **104**, May 2014, 5.
- [34] Sher, I., “Credibility and Determinism in a Game of Persuasion,” *Games and Economic Behavior*, **71**, March 2011, 409–419.
- [35] Shin, H. S., “Disclosures and Asset Returns,” *Econometrica*, **71**, January 2003, 105–133.
- [36] Spence, M., “Job Market Signaling,” *Quarterly Journal of Economics*, **87**, August 1973, 355–374.
- [37] Sugaya, T., and A. Wolitzky, “The Revelation Principle in Multistage Games,” MIT working paper, 2018.