# Mechanism Design with Limited Commitment[*]

Laura Doval[†]        Vasiliki Skreta[‡]

First version: April 8, 2018. This version: November 22, 2018

## Abstract

We develop a tool akin to the revelation principle for mechanism design with limited commitment. We identify a *canonical* class of mechanisms rich enough to replicate the payoffs of any equilibrium in a mechanism-selection game between an uninformed designer and a privately informed agent. A cornerstone of our methodology is the idea that a mechanism should encode not only the rules that determine the allocation, but also the information the designer obtains from the interaction with the agent. Therefore, how much the designer learns, which is the key tension in design with limited commitment, becomes an explicit part of the design. We show how this insight can be used to transform the designer's problem into a constrained optimization one: To the usual truthtelling and participation constraints, one must add the designer's sequential rationality constraint.

KEYWORDS: mechanism design, limited commitment, revelation principle, information design

JEL CLASSIFICATION: D84, D86

## 1  Introduction

The standard assumption in dynamic mechanism design is that the designer can commit to long-term contracts. This assumption is useful: It allows us to characterize the best possible payoff for the designer in the presence of adverse selection and/or moral hazard, and it is applicable in many settings. Often, however, this assumption is done for technical convenience. Indeed, when the designer can commit to long-term contracts, the mechanism-selection problem can be reduced to a constrained optimization problem thanks to the *revelation principle*.[1] However, as the literature starting with Laffont and Tirole (1987, 1988) shows, when the designer can only commit to short-term contracts, the tractability afforded by the revelation principle is lost. Indeed, mechanism design problems with limited commitment are difficult to analyze without imposing auxiliary assumptions either on the class of contracts the designer can choose from, as in Gerardi and Maestri (2018) and Strulovici (2017), or on the length of the horizon, as in Skreta (2006, 2015).

This paper provides a "revelation principle" for dynamic mechanism-selection games in which the designer can only commit to short-term contracts. We study a game between an uninformed designer and an informed agent with persistent private information. Although the designer can commit within each period to the terms of the interaction–the current mechanism–he cannot commit to the terms the agent faces later on, namely, the mechanisms that are chosen in the continuation game. First, we show there is a class of mechanisms that is sufficient to replicate all equilibrium payoffs of the mechanism-selection game. Second, we show how this insight can be used to transform the designer's problem into a constrained optimization one: To the usual truthtelling and participation constraints, one must add the designer's sequential rationality constraint.

The starting point of our analysis is the class of mechanisms we allow the designer to select from. Following Myerson (1982) and Bester and Strausz (2007), we consider mechanisms defined by a *communication device* and an *allocation rule*

---

[1]The "revelation principle" denotes a class of results in mechanism design; see Gibbard (1973), Myerson (1979), and Dasgupta et al. (1979).
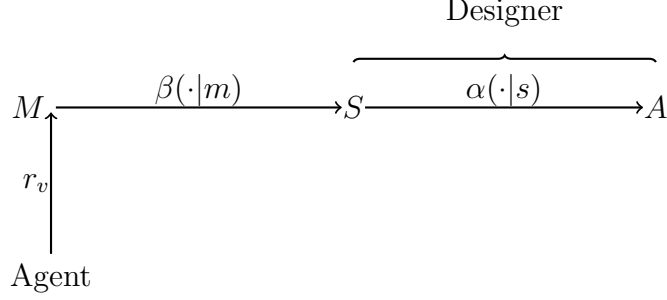
as illustrated in Figure 1:[2]



Figure 1: Mechanisms: communication device, $(M, \beta, S)$, and allocation rule, $\alpha$

Having observed her private information (her type, $v \in V$), the agent privately reports an *input* message, $m \in M$, into the mechanism; this then determines the distribution, $\beta(\cdot|m)$, from which an *output* message, $s \in S$, is drawn. In turn, the output message determines the distribution, $\alpha(\cdot|s)$, from which the allocation is drawn. The output message and the allocation are publicly observable: They constitute the contractible parts of the mechanism.

When the designer has commitment power, the revelation principle implies that, without loss of generality, we can restrict attention to mechanisms satisfying the following three properties: (i) $M = V$, (ii) $M = S$, and (iii) $\beta$ is "invertible." By $\beta$ being "invertible," we mean the designer learns the input message by observing the output message; in this case, the designer learns the agent's type report upon observing the output message. Moreover, the revelation principle implies we can restrict attention to equilibria in which the agent reports her type truthfully, which means the designer not only learns the agent's type report upon observing the output message, but also he learns the agent's true type.

It is then clear why restricting attention to mechanisms that satisfy properties (i)-(iii) and truthtelling equilibria is with loss of generality under limited commitment: Upon observing the output message, the designer learns the agent's type

---
[2]Myerson (1982) allows the designer to choose, as a function of the input message $m$, any joint distribution over $S \times A$. It is a consequence of Theorem 3.1 that the mechanisms in Figure 1 are without loss of generality; see Appendix II for a proof. Since the formulation of a mechanism in Figure 1 allows us to highlight the role of the communication device separately from that of the allocation, we opt for this formulation for pedagogical purposes.

report and hence her type. Then the agent may have an incentive to misreport if the designer cannot commit not to react to this information. This is precisely the intuition behind the main result in Bester and Strausz (2001), which is the first paper to provide a general analysis of optimal mechanism design with limited commitment. The authors restrict attention to mechanisms in which the cardinality of the set of input and output messages is the same and $\beta$ is "invertible." They show that to sustain payoffs in the Pareto frontier, mechanisms in which input messages are type reports are without loss of generality. However, focusing on truthtelling equilibria is with loss of generality. In a follow-up paper, Bester and Strausz (2007) lift the restriction on the class of mechanisms (i.e., (ii) and (iii) above) and show in a one-period model that focusing on mechanisms in which input messages are type reports and truthtelling equilibria is without loss of generality. The authors, however, do not characterize the output messages. It is also not clear whether taking input messages to be type reports is without loss when the designer and the agent interact repeatedly (see the discussion after Theorem 3.1).

The main contribution of this paper is to show that, under limited commitment, taking the set of output messages to be the set of posterior beliefs of the designer about the agent's type, that is, $S = \Delta(V)$, is without loss of generality. Theorem 3.1 shows that in a general mechanism-selection game between an uninformed designer and an informed agent introduced in Section 2, any equilibrium payoff can be replicated by an equilibrium in which (a) the designer uses mechanisms in which input messages are type reports and output messages are beliefs, (b) the agent always participates in the mechanism, and (c) input and output messages have a *literal* meaning: The agent reports her type truthfully, and if the mechanism outputs a given posterior, this posterior coincides with the belief the designer holds about the agent's type given the agent's strategy and the mechanism. Given that any equilibrium payoff can be replicated by mechanisms in which input messages are type reports and output messages are beliefs about the agent's type, we call this class of mechanisms *canonical*.

Theorem 3.1 implies that when the designer is subject to sequential rationality constraints, the mechanism serves a dual role within a period. On the one hand, it determines the allocation for that period. On the other hand, it determines the in-

formation about the agent that is carried forward in the interaction. An advantage of the language of posterior beliefs is that it avoids potential infinite-regress problems. Indeed, in a finite horizon problem, an alternative set of output messages could be a recommendation for an allocation today and a sequence of allocations from tomorrow on.[3] In the final period, the revelation principle in Myerson (1982) pins down the implementable allocations. Therefore, the recommended allocations can be determined via backward induction. This idea cannot be carried to an infinite horizon setting: These sets of output messages would necessarily have to make reference to the continuation mechanisms, which are themselves defined by a set of output messages.

Another contribution of our analysis is to show that to characterize equilibrium payoffs of the game between the designer and the agent, it suffices to consider a simpler game, denoted the *canonical game*. We record this result in Proposition 3.1. In the canonical game–studied in Section 3.2–the designer is restricted to offer mechanisms in which input messages are type reports and output messages are beliefs over the agent's type. Theorem 3.1 (trivially) implies an equilibrium outcome of the canonical game can be achieved by strategy profiles in which the principal employs mechanisms that induce the agent to truthfully report his type and to always participate. However, the principal has *fewer deviations* in the canonical game and an equilibrium strategy may not be an equilibrium if the principal can deviate to *any* mechanism, as he can in the mechanism-selection game. One may then wonder whether analyzing the canonical game gives, unintentionally, some commitment power to the principal.

Proposition 3.1 shows this is not the case: Leveraging the construction used to establish Theorem 3.1, we show that, without loss of generality, the *best deviation* in the mechanism-selection game is equivalent to a deviation to a canonical mechanism that induces the agent to report truthfully and to participate with probability one. In a finite horizon setting, Proposition 3.1 justifies writing the designer's problem as a sequence of maximization problems over canonical mechanisms subject to the agent's participation and incentive compatibility constraints and the designer's sequential rationality constraints.

---

[3]See Section 5.1 for a formal discussion of the approach and its potential issues.

Section 4 illustrates the methodology for the case of transferable utility and preferences that satisfy increasing differences in distributions. The resulting program allows us to highlight the connection between our problem and the literature on information design; after all, the designer can be thought of as a sender who designs the information structure for a receiver, who happens to be his future self. However, there are differences. In our setting, the first-period principal (the sender in Kamenica and Gentzkow (2011)) also takes an action for each posterior he induces. In addition, the first-period principal's objective function depends on the prior as well as the posterior, whereas in Kamenica and Gentzkow (2011), it only depends on the posterior. Finally, the first-period principal cannot implement any Bayes' plausible distribution over posteriors, but only those that satisfy the incentive compatibility and participation constraints of the agent.

An important difference between the mechanisms used by Hart and Tirole (1988), Laffont and Tirole (1988), Freixas et al. (1985), and Bester and Strausz (2001) and the ones considered here is that whereas in the former papers, the principal observes the agent's choice out of a menu of *contracts*, here, the agent's input into the communication device is not observed. Under the assumptions of Section 4, Proposition 5.2 in Section 5.2 characterizes the mechanisms (i.e., the communication device and allocation pairs) that can be implemented with the agent making a choice out of a menu. The result is useful for the following reasons. First, by checking whether the solution to the program studied in Section 4 satisfies the conditions in Proposition 5.2, we can understand whether the modeling of a mechanism as a menu of contracts in the aforementioned works is without loss. Second, when the solution to the program does satisfy the conditions, it allows the analyst to propose a "simple" implementation of the optimum.

The paper contributes to the literature on mechanism design with limited commitment, referenced throughout the introduction.[4] A large literature studies the effect of limited commitment within a specific class of "mechanisms": The papers in the durable-good monopolist literature (Bulow (1982); Gul et al. (1986); Stokey

---

[4]A designer's lack of commitment can take various forms, not considered in this paper, but that have been studied in other papers. See, for instance, McAdams and Schwarz (2007), Vartiainen (2013), and Akbarpour and Li (2018), in which the designer cannot commit even to the obey the rules of the current mechanism.

(1981)) study price dynamics and establish (under some conditions) Coase's conjecture whereby a monopolist essentially loses all profits if he lacks commitment. In an analogous vein, Burguet and Sakovics (1996), McAfee and Vincent (1997), Caillaud and Mezzetti (2004), and Liu et al. (2018) study equilibrium reserve-price dynamics without commitment in different setups. The common thread is, again, that the seller's inability to commit reduces monopoly profits.

Mechanism-selection in a dynamic environment with limited commitment is considered in Deb and Said (2015). The authors study a model of sequential screening, in which new buyers arrive over time. Like in Skreta (2006) and Skreta (2015), Deb and Said (2015) consider general mechanisms but a finitely long interaction. Infinitely long contract-selection games are studied in Strulovici (2017) and Gerardi and Maestri (2018). The former studies renegotiation and finds that equilibrium allocations become efficient as the parties become arbitrarily patient. In Gerardi and Maestri (2018), however, the limit allocation is inefficient whenever firing the agent–what the authors refer to "firing allocation"– is not a solution when there is commitment.

By highlighting the role that the designer's beliefs about the agent play in mechanism design with limited commitment, our paper also relates to Lipnowski and Ravid (2017) and Best and Quigley (2017), who study models of direct communication between an informed sender and an uninformed receiver.[5] Lipnowski and Ravid (2017) show how the posterior approach of Kamenica and Gentzkow (2011) can be used to characterize equilibrium outcomes, and study their properties in the cheap talk model of Crawford and Sobel (1982) (the leading model of communication without commitment), when the sender's preferences do not depend on the state of the world.[6] Finally, given that the search for the best equilibrium often reduces to solving a constrained information design problem we relate to, among others, Le Treust and Tomala (2017), Georgiadis and Szentes (2018), and Boleslavsky and Kim (2018).

---

[5]Salamanca (2016) studies mediated communication in Kamenica and Gentzkow (2011).

[6]Golosov and Iovino (2016) study a social insurance model with a continuum of agents, where private information is not persistent across stages. They leverage the resulting repeated-game structure to solve for the best equilibrium.

The rest of the paper is organized as follows. Section 2 describes the model and notation. Section 2.1 analyzes a simple version of the model in Skreta (2006); it allows us to introduce the main ideas of the paper in a simple and well-known setting. Section 2.2 discusses the modeling assumptions. Section 3 introduces the main theorem and provides a sketch of the proof. Section 4 specializes the results to the two-period model of Bester and Strausz (2007) with transferable utility and single-crossing preferences. We compare the solution of the 'relaxed' problem to the information design model of Kamenica and Gentzkow (2011). Section 5.1 discusses using recommendations as output messages. Section 5.2 studies implementation when the principal observes the agent's choice. Section 5.3 discusses an example with multiple agents. All proofs are relegated to the Appendix. The supplementary material (Sections I-IV) contains omitted proofs and extensions discussed throughout the main text.

## 2 MODEL

**Primitives** There are two players: a principal (he) and an agent (she). They interact over $T \leqslant \infty$ periods. Before the game starts, the agent observes her type, $v \in V$. $V$ is any finite set; however, the main insights extend to the case in which $V$ is a Polish space (see Appendix IV). Each period, as a result of the interaction between the principal and the agent, an allocation $a \in A$ is determined. Assume $A$ is a compact (possibly finite) space.

Given a sequence of allocations $a^t = (a_0, a_1, ...., a_t)$, the principal can only choose $a_{t+1} \in \mathcal{A}(a^t)$. That is, there is a correspondence $\mathcal{A} : \bigcup_{n=0}^{T} A^n \mapsto A$ such that for $t \in \mathbb{N}, a^t \in A^t$, $\mathcal{A}(a^t)$ describes the set of allocations the principal can offer given the allocations he has offered in the past. Assume $\mathcal{A}$ is compact-valued and there exists an allocation $a^* \in A$ such that $a^*$ is always available.[7]

Payoffs are defined as follows. For the principal, assume there exists a function,

---

[7]We later use allocation $a^*$ to model the agent's participation decision within each period: If the agent prefers not to participate, allocation $a^*$ is implemented automatically. For instance, in a trade model such as the one in Section 2.1, $a^*$ corresponds to no trade and no transfers. The constraint correspondence $\mathcal{A}$ also allows us to capture that the agent can walk away from the mechanism as in Gerardi and Maestri (2018): We could specify that the first time $a^*$ is implemented, then this allocation is the only one available thereafter.

$W : A^T \times V \mapsto \mathbb{R}$ such that his payoff from allocation $a \in A^T$ when the agent's type is $v$ is given by $W(a, v)$. Similarly for the agent, when her type is $v$, her payoff from allocation $a \in A^T$ is given by $U(a, v)$.

**Mechanisms:** In each period, the principal offers the agent a mechanism, $\mathbf{M}_t = \langle (M^{\mathbf{M}_t}, \beta^{\mathbf{M}_t}, S^{\mathbf{M}_t}), \alpha^{\mathbf{M}_t} \rangle$, which consists of a *communication device*, $(M^{\mathbf{M}_t}, \beta^{\mathbf{M}_t}, S^{\mathbf{M}_t})$, and an *allocation rule*, $\alpha^{\mathbf{M}_t}$, where

$$\beta^{\mathbf{M}_t} : M^{\mathbf{M}_t} \mapsto \Delta^*(S)$$
$$\alpha^{\mathbf{M}_t} : S^{\mathbf{M}_t} \mapsto \Delta^*(\mathcal{A}),$$

and where $\Delta^*(C)$ denotes the set of distributions on $C$ with *finite* support. We endow the principal with a collection $(M_i, S_i)_{i \in \mathcal{I}}$ of input and output message sets in which each $M_i$ is finite, $|V| \leqslant |M_i|$, and $\Delta(M_i) \subseteq S_i$.[8] Moreover, we assume $(V, \Delta(V))$ is an element in that collection. Denote by $\mathcal{M}$ the set of all mechanisms with message sets $(M_i, S_i)_{i \in \mathcal{I}}$. A mechanism is *canonical* if $(V, \Delta(V))$ are its sets of input and output messages. Let $\mathcal{M}^C$ denote the set of canonical mechanisms and let $\mathbf{M}_t^C$ denote an element in that set.

Three remarks are in order. First, the restriction that $M_i$ has at least as many messages as types is without loss of generality. The principal can always replicate a mechanism with a smaller set of input messages by using a larger set of input messages.[9] Second, we restrict the principal to design $\beta^{\mathbf{M}_t}$ and $\alpha^{\mathbf{M}_t}$ to be distributions with finite support, thus allowing us to focus on the novel conceptual features of the environment, as opposed to dealing with measure-theoretic complications. To replicate any equilibrium of the game when the principal selects distributions with finite support using canonical mechanisms, we find the principal only needs to use distributions with finite support. This last observation, of course, would not be true if the set of types were not finite.[10],[11] Finally, we re-

---

[8]Technically, we only need that $S_i$ contains an image of $\Delta(M_i)$.

[9]To see this, suppose the principal would rather use a mechanism, $\mathbf{M}_t'$, with a message space $M^{\mathbf{M}_t'}$ with cardinality strictly less than $|V|$. Then he can choose a mechanism $\mathbf{M}_t$ with $M^{\mathbf{M}_t} = V$, choose $\beta$ to coincide with $\beta^{\mathbf{M}_t'}$ on the first $|M^{\mathbf{M}_t'}|$ messages, and have $\beta^{\mathbf{M}_t}$ coincide with $\beta^{\mathbf{M}_t'}(\cdot|m_1')$ for all remaining messages.

[10]Appendix IV extends our result to the case in which $V$ is a compact and metrizable space.

[11]We conjecture, however, that the restriction to distributions with finite support is without

strict the principal to choose input and output messages within the set $(M_i, S_i)_{i \in \mathcal{I}}$. This allows us to have a well-defined set of deviations for the principal, avoiding set-theoretic issues related to self-referential sets. The analysis that follows shows that the choice of the collection plays no further role in the analysis.

**Timing:** In each period $t$,

- The principal and the agent observe a draw from a correlating device $\omega \sim U[0, 1]$.

- The principal offers the agent a mechanism $\mathbf{M}_t$.

- The agent observes the mechanism and decides whether to participate ($p = 1$) or not ($p = 0$). If she does not participate, $a^*$ is implemented and the game proceeds to $t + 1$.

- If she participates, she privately submits a report $m \in M^{\mathbf{M}_t}$.

- $s \in S^{\mathbf{M}_t}$ is drawn according to $\beta^{\mathbf{M}_t}(\cdot|m)$, which is publicly observed.

- $a \in A$ is drawn according to $\alpha^{\mathbf{M}_t}(\cdot|s)$, which is publicly observed.

This defines an extensive form game, which we dub the *mechanism-selection game*. If, instead, the principal can only choose mechanisms in $\mathcal{M}^C$, we denote it as the *canonical game*.

Public histories in this game are[12]

$$h^t = (\omega_0, \mathbf{M}_0, p_0, s_0, a_0, \ldots, \omega_{t-1}, \mathbf{M}_{t-1}, p_{t-1}, s_{t-1}, a_{t-1}, \omega_t),$$

where $p_r \in \{0, 1\}$ denotes the agent's participation with the restriction that $p_r = 0 \Rightarrow s_r = \varnothing, a_r = a^*$. Public histories capture what the principal knows through period $t$. Let $H^t$ denote the set of all period $t$ public histories. A strategy for the principal is then given by $\Gamma_t : H^t \mapsto \Delta(\mathcal{M})$.

---

loss of generality when the set of types is finite.

[12]The restriction that the support of $\beta^{\mathbf{M}_t}(\cdot|m)$ is finite for $m \in M^{\mathbf{M}_t}$, together with the finiteness of $M^{\mathbf{M}_t}$, imply that there are output messages $s \in S^{\mathbf{M}_t}$ that can never arise. Thus, we can remove from the tree all the histories that are consistent with mechanism $\mathbf{M}_t$ being offered and $s \in S^{\mathbf{M}_t}$ such that $\sum_{m \in M^{\mathbf{M}_t}} \beta^{\mathbf{M}_t}(s|m) = 0$, without affecting the equilibrium set. However, for tractability, we do not make this distinction in our notation.

A history for the agent consists of the *public* history of the game together with the agent's inputs into the mechanism (henceforth, the agent history) and her private information. Formally, an agent history is an element

$$h_A^t = (\omega_0, \mathbf{M}_0, m_0, p_0, s_0, a_0, \ldots, \omega_{t-1}, \mathbf{M}_{t-1}, p_{t-1}, m_{t-1}, s_{t-1}, a_{t-1}, \omega_t).$$

Given a public history $h^t$, let $H_A^t(h^t)$ denote the set of agent histories consistent with $h^t$. The agent also knows her type, and hence a history through period $t$ is an element of $\{v\} \times H_A^t$ when her type is $v$. The agent's participation strategy is $\pi_v : H_A^t \times \mathcal{M}_t \mapsto [0,1]$. Conditional on participating in the mechanism $\mathbf{M}_t$, her reporting strategy is a distribution $r_v(h_A^t, \mathbf{M}_t, 1) \in \Delta(M^{\mathbf{M}_t})$ for each of her types $v$ and each $h_A^t \in H_A^t$.

A belief for the principal at the beginning of time $t$, history $h^t$, is a distribution $\mu(h^t) \in \Delta(V \times H_A^t(h^t))$, where $H_A^t(h^t)$ is the set of agent histories that are consistent with the public history $h^t$, which is observed by the principal. The principal is thus uncertain both about the agent's payoff-relevant type, $v$, and her payoff-irrelevant type, $h_A^t$.

Our focus is on studying the equilibria of the mechanism-selection and canonical games. By equilibrium, we mean Perfect Bayesian equilibrium (henceforth, PBE), defined as follows:

**Definition 2.1.** A *Perfect Bayesian Equilibrium* is a tuple $\langle \Gamma^*, (\pi_v^*, r_v^*)_{v \in V}, \mu^* \rangle$ such that for each $h^t$ the following hold:

1. Given $\mu_t^*(h^t)$, $\Gamma_t^*(h^t)$ is sequentially rational given $(\pi_v^*, r_v^*)_{v \in V}$,

2. Given $\Gamma^*(h^t)$, $\pi_v^*(h_A^t, \cdot), r_v^*(h_A^t, \cdot, 1)$ are sequentially rational for all $h_A^t \in H_A^t(h^t)$,

3. $\mu^*(h^t)$ is derived via Bayes' rule whenever possible.

Implicit in the definition of PBE is the assumption that the principal does not update his beliefs about the agent following a deviation by the principal. That is, we assume beliefs are *pre-consistent* (see Hendon et al. (1996)).

**Remark 2.1.** [Belief updating depends only on the realized output message] Fix

a history $h^t$. Given $\mu \in \Delta(V \times H_A^t(h^t))$ and a mechanism $M^{\mathbf{M}_t}$, Bayesian updating depends on the agent's strategy and the communication device, but not on the allocation rule. To see this, suppose the agent participates with positive probability in the mechanism, the output message is $s \in S^{\mathbf{M}_t}$ and allocation $a$ is observed; then the principal's belief about the agent being at history $(v, h_A^t, \mathbf{M}_t, 1, s, a)$ is given by

$$\frac{\mu(h^t)(v, h_A^t)\pi_v^*(h_A^t, \mathbf{M}_t) \sum_{m \in M^{\mathbf{M}_t}} r_v^*(h_A^t, \mathbf{M}_t, 1)(m)\beta^{\mathbf{M}_t}(s|m)\alpha^{\mathbf{M}_t}(a|s)}{\sum_{\tilde{v}, \widetilde{h_A^t}} \mu(\tilde{v}, \widetilde{h_A^t})\pi_{\tilde{v}}^*(\widetilde{h_A^t}, \mathbf{M}_t) \sum_{\tilde{m} \in M^{\mathbf{M}_t}} r_{\tilde{v}}^*(\widetilde{h_A^t}, \mathbf{M}_t, 1)(\tilde{m})\beta^{\mathbf{M}_t}(s|\tilde{m})\alpha^{\mathbf{M}_t}(a|s)},$$

and all the terms concerning $\alpha^{\mathbf{M}_t}(a|s)$ drop out.

## 2.1  Example: two-period sale of a durable good.

To fix ideas, consider the following example. The principal is a seller who owns one unit of a durable good and assigns value 0 to it. The agent is a buyer whose valuation for the good is her private information. The buyer's valuation can take two values, $v \in \{v_L, v_H\}$, where $v_H - v_L > 0$. The seller's belief that $v = v_H$ is $\mu_1$. An allocation is a pair $(q, t) \in \{0, 1\} \times \mathbb{R}$, where $q$ indicates whether a sale occurs ($q = 1$) or not ($q = 0$), and $t$ is a transfer from the buyer to the seller. Utilities are quasilinear so that the buyer's utility is $u(q, t; v) = vq - t$ and the seller's is $w(q, t; v) = t$. Both players share a common discount factor $\delta \in (0, 1)$.

The timing is as follows: in each period $t \in \{1, 2\}$

- The seller chooses a mechanism.

- The buyer observes the mechanism and decides whether to participate.

  – If she does not participate, the good is not sold and no payments are made; if $t = 1$, we move to period 2.

  – If she participates, the mechanism determines the allocation.

- If the good is not sold and $t = 1$, move on to $t = 2$.

Because the horizon is finite, we can solve the game by backward induction. Then let $t = 2$ and denote by $\mu_2$ the seller's posterior belief that $v = v_H$. In $t = 2$,

the seller has full commitment and the solution is routine. The seller posts a price equal to $v_L$ when $\mu_2 < v_L/v_H \equiv \overline{\mu}$, a price equal to $v_H$ when $\mu_2 > \overline{\mu}$, and at $\mu_2 = \overline{\mu}$, then the seller is indifferent between the two prices. Thus, the seller's revenue as a function of $\mu_2$ is given by

$$R_2(\mu_2) = \begin{cases} v_L & \text{if } \mu_2 \leqslant \overline{\mu} \\ \mu_2 v_H & \text{otherwise} \end{cases} = \begin{cases} \mu_2 v_H + (1 - \mu_2)\hat{v}_L(\mu_2) & \text{if } \mu_2 \leqslant \overline{\mu} \\ \mu_2 v_H & \text{otherwise} \end{cases},$$

where $\hat{v}_L(\mu_2) = v_L - (\mu_2/(1 - \mu_2))(v_H - v_L)$ and the equality follows from noting that when the price is $v_L$, the seller leaves rents $v_H - v_L$ with probability $\mu_2$ to the high type.

We now turn to period 1. Recall that $\mu_1$ denotes the probability that the buyer's valuation is $v_H$. Consistent with the mechanism-selection game introduced in Section 2, we allow the seller to offer the buyer a mechanism that consists of a communication device $(M_1, \beta_1, S_1)$ and an allocation rule $\alpha : S_1 \mapsto \Delta(\{0, 1\} \times \mathbb{R})$. The assumption of quasilinearity implies that, without loss of generality, the seller does not randomize on the transfers, so that $\alpha_1(q, t|s_1) = q(s_1) \times \mathbb{1}[t = t(s_1)]$.

Theorem 3.1 shows that, without loss of generality, input messages are type reports, $M_1 = V$, and output messages are the seller's beliefs about the buyer's valuation, $S_1 = \Delta(V)$. We now provide intuition for this in the context of the example.

1. To see that $M_1 = V$, note that $\beta_1$ together with the agent's reporting strategy induces another distribution on $S_1$,

$$\sum_{m_1 \in M_1} \beta_1(s_1|m_1) r_v(m_1) \equiv \beta^*(s_1|v).$$

   If the seller offers $\langle (V, \beta^*, S_1), \alpha_1 \rangle$ to the buyer, then the buyer tells the truth (see also Bester and Strausz (2007)).

2. To see why $S_1$ can be taken to be $\Delta(V)$, note that upon the realization of $s_1$, two things happen. First, the allocation $\alpha(s_1)$ is determined. Second, if the

13

allocation is no trade, $s_1$ is used to update the principal's beliefs as follows:

$$\mu_2(v = v_H|s_1)\left(\sum_{v\in V}\mu_1(v)\beta^*(s_1|v)\right) = \mu_1(v_H)\beta^*(s_1|v_H),$$

where we have already used that $M_1 = V$ and the buyer reports truthfully. Given the belief induced by $s_1$, we know what happens in period 2; there is no use for $s_1$ beyond that. Thus, we can take $S_1 = \Delta(V)$. Thus, we write $\beta(\mu_2|v), q(\mu_2), t(\mu_2)$ instead of $\beta(s_1|v), q(s_1), t(s_1)$ thereafter.

With these observations, we can describe the seller's optimal mechanism in period 1 via the following program:

$$R_1(\mu_1) \equiv \max_{\beta,q,t} \sum_{\mu_2\in\Delta(V)}\left(\sum_{v\in V}\mu_1(v)\beta(\mu_2|v)\right)[t(\mu_2) + (1 - q(\mu_2))\delta R_2(\mu_2)]$$

subject to for all $v \in \{v_L, v_H\}$:

$$\sum_{\mu_2\in\Delta(V)}\beta(\mu_2|v)(vq(\mu_2) - t(\mu_2) + \delta(1 - q(\mu_2))u_B(\mu_2, v)) \geqslant 0 \ (\text{PC}_v)$$

$$\sum_{\mu_2\in\Delta(V)}(\beta(\mu_2|v) - \beta(\mu_2|v'))(vq(\mu_2) - t(\mu_2) + \delta(1 - q(\mu_2))u_B(\mu_2; v)) \geqslant 0 \ (\text{IC}_{v,v'})$$

$$\mu_2(v_H)\left(\sum_{v\in V}\mu_1(v)\beta(\mu_2|v)\right) = \mu_1(v_H)\beta(\mu_2|v_H) \ (\text{BC}_{\mu_2}).$$

That is, the seller chooses $\beta, q, t$ to maximize his profit subject to the agent's participation and incentive compatibility constraint and a Bayesian consistency constraint. The latter says that when the mechanism outputs $\mu_2$, then $\mu_2$ is the belief that obtains via Bayesian updating. The buyer's participation and incentive compatibility constraints take into account her continuation values, denoted by $u_B(\mu_2, v)$: For low values of $\mu_2$, the high type is served at a low price in period 2.[13]

---

[13]Implicit in the buyer's participation constraint is that, if she does not participate in $t = 1$, the seller has belief $\mu_2 \geqslant \overline{\mu}$ and then sets a price of $v_H$ in $t = 2$. Thus, both types of the agent earn a payoff of 0 in case they do not participate in $t = 1$. Given Theorem 3.1, this is without loss of optimality. Indeed, Theorem 3.1 shows that it is without loss of generality to have the agent participate with probability 1. Hence, not participating of the mechanism becomes an off-path event and beliefs are not pinned down by Bayes' rule in this case.

As usual, we can show that $PC_{v_L}$ and $IC_{v_H,v_L}$ bind, and these two constraints imply the others. Therefore, we can use them to replace the transfers in the seller's objective to obtain

$$R_1(\mu_1) \equiv \max_{\tau,q} \sum_{\tau \in \Delta(V)} \tau(\mu_2) \left[ q(\mu_2)(\mu_2 v_H + (1 - \mu_2)\hat{v}_L(\mu_1)) + \delta(1 - q(\mu_2))R_2(\mu_2; \mu_1) \right]$$

(1)

$$\text{s.t.} \quad \sum_{\mu_2 \in \Delta(V)} \tau(\mu_2)\mu_2 = \mu_1,$$

where $\tau(\mu_2) = \sum_{v \in V} \mu_1(v)\beta(\mu_2|v)$ is the probability that $\mu_2$ is the induced posterior and

$$\delta R_2(\mu_2; \mu_1) = \begin{cases} \delta(\mu_2 v_H + (1 - \mu_2)\hat{v}_L(\mu_1)) & \text{if } \mu_2 < \overline{\mu} \\ \delta\mu_2 v_H & \text{if } \mu_2 > \overline{\mu} \end{cases}$$

(2)

is an adjusted version of the seller's period 2 revenue. We now explain equations (1) and (2) in detail. Equation (1) shows that the seller's period 1 problem can be solved by finding (i) a trade probability for each posterior and (ii) a distribution over posteriors that averages out to the prior. Given a posterior $\mu_2$, the trade probability, $q(\mu_2)$, is chosen to maximize a version of the virtual surplus, familiar from mechanism design with commitment. To see this, note that in equation (1), the probability of each type is evaluated using the posterior $\mu_2$, but the virtual value for the low type is computed using the prior $\mu_1$. This reflects that the seller in period 1 assigns probability $\mu_1$ to $v = v_H$, and $\mu_1$ is the rate at which he pays rents to the high type. Similarly, $\delta R_2(\mu_2; \mu_1)$ adjusts the revenues in period 2 by the rents the period 1 seller must leave to the buyer: In effect, should the period 1 seller induce $\mu_2 < \overline{\mu}$, the buyer obtains a rent of $\delta(v_H - v_L)$, which the seller in period 1 has to take into account.

Note that equation (2) does not specify what the seller's payoff is when $\mu_2 = \overline{\mu}$. This, in fact, depends on the prior $\mu_1$: When $\mu_2 = \overline{\mu}$, the period 2 seller is indifferent between prices $v_H$ and $v_L$. The period 1 seller, however, is not indifferent; this fact is illustrated in Figures 2 and 3 below:
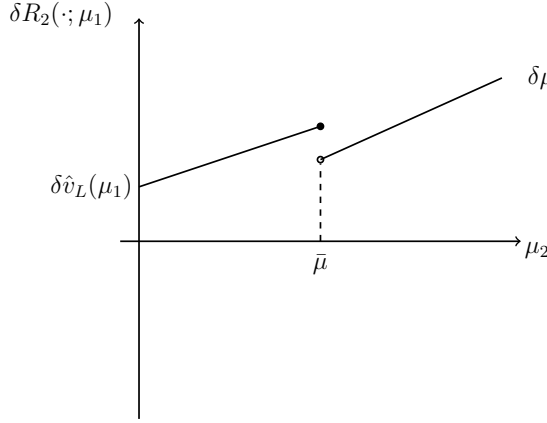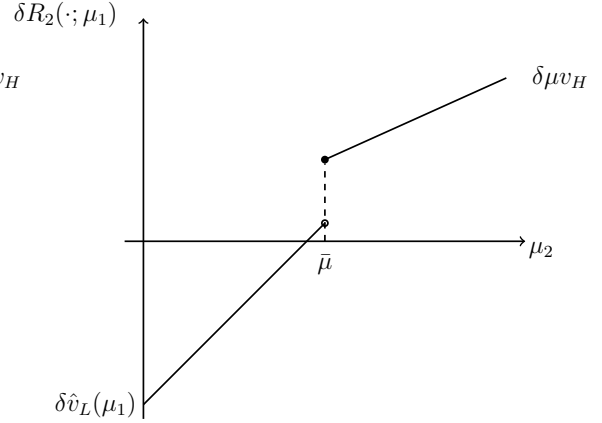
15

Figure 2: $\mu_1 \leqslant \bar{\mu}$



Figure 3: $\bar{\mu} < \mu_1$

If the seller's prior is such that he would sell to the low valuation buyer today ($\mu_1 \leqslant \bar{\mu}$), then he would rather have the period 2 seller also serve the low valuation type when indifferent in period 2, as illustrated in Figure 2. However, if the period 1 seller would prefer to exclude the high valuation buyer when her valuation is low, then he would prefer the low valuation buyer to be excluded in period 2 as well when $\mu_2 = \bar{\mu}$, as illustrated in Figure 3.

In what follows, we solve the seller's problem for the case in which $\mu_1 > \bar{\mu}$.[14] Because the seller can choose $q(\mu_2)$ for each $\mu_2$, the best he can do is choose it to pointwise maximize the objective function in equation (1), as illustrated in Figures 4 and 5 below:

---

[14]The case in which $\mu_1 < \bar{\mu}$ is immediate: The seller can achieve the commitment solution by selling to both types of the buyer in period 1.
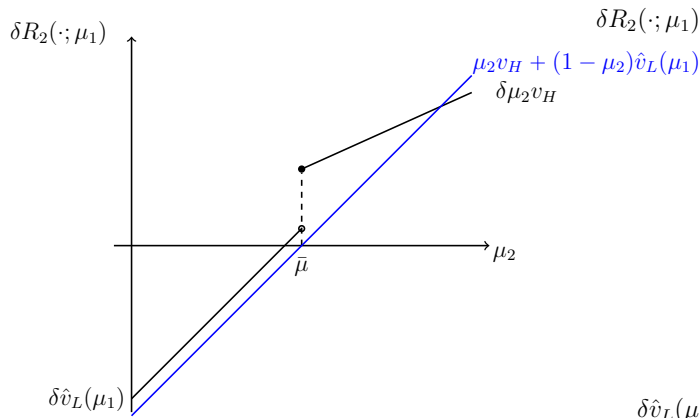
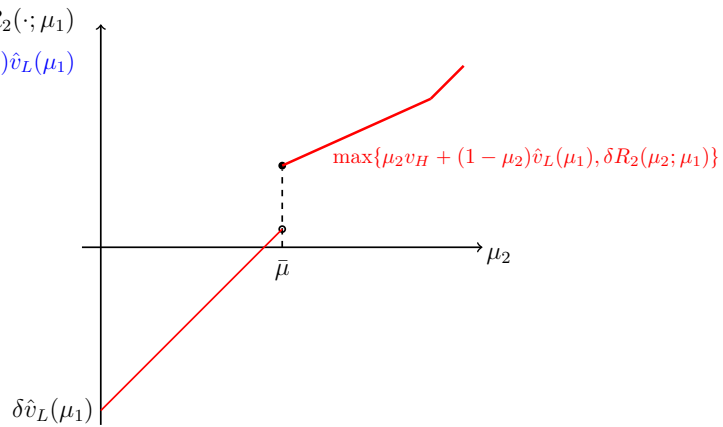Figure 4: Value of setting $q(\mu_2) = 0$ (black) and $q(\mu_2) = 1$ (blue)



Figure 5: Pointwise maximum of the blue and black lines in Figure 4

This reduces the principal's problem to that of finding a distribution over posteriors to solve:

$$\max_{\tau} \sum_{\mu_2 \in \Delta(V)} \tau(\mu_2) \max\{\mu_2 v_H + (1 - \mu_2)\hat{v}_L(\mu_1), \delta R(\mu_2; \mu_1)\},$$

subject to the constraint that the distribution must average to the prior. Under some parameter values, the solution is as depicted in Figure 6 below:[15]

---

[15]When $\overline{\mu} < \mu_1$, there are two possible solutions depending on the parameter values. When $\mu_1 > \overline{\mu}$ is high enough, we obtain the solution depicted in Figure 6 and described in the main text. For lower values of $\mu_1 > \overline{\mu}$, we obtain the solution familiar to the literature on the ratchet effect (see Hart and Tirole (1988)). In this case, the seller sets a price of $v_L$ in period 2, and a price of $v_H - \delta \Delta v$ in period 1; the buyer buys in period 1 when $v = v_H$ and in period 2 when $v = v_L$.

The plot shows axes labeled $\delta R_2(\cdot;\mu_1)$ (vertical) and $\mu_2$ (horizontal), with curve label $\max\{\mu_2 v_H + (1-\mu_2)\hat{v}_L(\mu_1), \delta R_2(\mu_2;\mu_1)\}$, horizontal axis marks $\overline{\mu}$ and $\mu_1$, and vertical intercept $\delta\hat{v}_L(\mu_1)$.
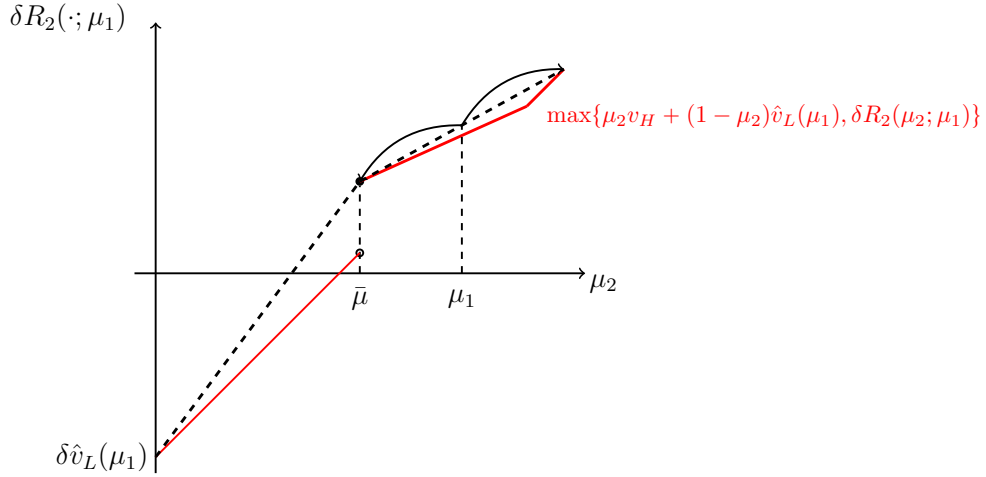
Figure 6: Optimal mechanism for period 1 seller: The black dashed line depicts the concavification of the function in Figure 5

The seller sets a price of $v_H$ in periods 1 and 2. The buyer does not buy when her value is low, whereas the buyer randomizes between buying today and tomorrow when her value is high. The randomization is such that when the seller sees no sale at the end of period 1, he attaches probability $\overline{\mu}$ to $v = v_H$.

The example highlights both how the language of type reports and posterior beliefs is enough to replicate what the principal can obtain from any other mechanism and also how useful this language is to solve mechanism design problems with limited commitment. Indeed, it allows us to reduce the problem of finding the best equilibrium for the principal to a constrained optimization problem. In Section 4, we return to the setting of transferable utility and preferences that satisfy increasing differences and show that the connection between our problem and information design extends beyond the example.

However, the example does not allow us to highlight some features of the model, which we discuss in the next section. The reader eager to see the results can skip to Section 3; however, the discussion may be useful to follow the proof sketch of the main theorem.

## 2.2 Discussion: Randomized allocations and public correlating device

We now discuss two aspects of the model that do not seem to play a role in the example, but are important in what follows: The principal is allowed to offer a randomization over allocations, and the principal and the agent have access to a public correlating device.

**Randomized allocations** There are two reasons for allowing the principal to choose randomized allocations. First, randomized allocations are necessary for the set of input messages to be the set of type reports; this is inherited from the revelation principle with commitment (see Strausz (2003)). To see this, consider the situation illustrated in Figure 7 below. The mechanism is simple: If the agent reports $m$, then the output message is $m$ and the allocation is $a$, whereas if she reports $m'$, the output is $m'$ and the allocation is $a'$. Assume that when her type is $v$, the agent sends $m$ and $m'$ with probability $p$ and $1-p$, respectively; thus, she obtains $a$ and $a'$ with probability $p$ and $1-p$, respectively.



Figure 7: Agent of type $v$ randomizes over $m$ and $m'$ generating a randomized allocation

If we restrict the principal to offer deterministic allocations, then he cannot replicate the agent's allocation just by asking for a truthful type report. However, if we allow the principal to offer a mechanism such that when the input is $v$ and the output is $v$, the allocation is a randomization between $a$ and $a'$, then he can replicate the allocation type $v$ obtains just by soliciting a type report.

Second, randomized allocations are necessary for the set of output messages to be the set of distributions over the agent's type. To see this, note that two different output messages, $s$ and $s'$, may be associated with two different allocations, $a$ and $a'$, but with the same posterior belief, as illustrated in Figure 8 below:

Figure 8: Two output messages, $s$ and $s'$, induce same posterior but different allocations
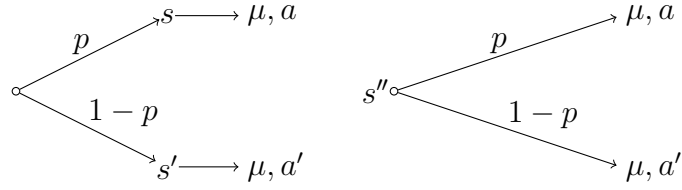
By allowing the principal to offer randomized allocations, we can collapse $s$ and $s'$ to one output message $s''$ associated to one posterior, $\mu$.

**Public correlating device** The correlating device is important for output messages to be the principal's posterior beliefs about the agent's type. Note that two output messages, $s$ and $s'$, may be associated with two different continuation equilibria, even if they induce the same allocation and posterior beliefs, as in Figure 9 below:



Figure 9: Two output messages, $s$ and $s'$, induce same posterior and allocations, but different continuation equilibria

The correlating device allows us to collapse $s$ and $s'$ into one output message (and hence, one posterior belief) and coordinate continuation play with the correlating device, akin to what is done in repeated games. This feature arises, somewhat trivially, in Section 2.1. In the example, for each posterior belief different from $\overline{\mu}$, there is a unique continuation equilibrium in period 2, and hence there is no need to select amongst continuation equilibria. However, when $\mu_2 = \overline{\mu}$, there are two continuation equilibria; when we allowed the first-period principal to select between them, we implicitly made use of a (trivial) correlating device.

## 3 Results

Section 3 presents the main results of the paper. Theorem 3.1 shows that any equilibrium payoff of the mechanism-selection game can be replicated by an equilibrium in which (a) the designer uses mechanisms in which input messages are type reports and output messages are beliefs, (b) the agent always participates in the mechanism, and (c) input and output messages have a *literal* meaning: The agent reports her type truthfully, and if the mechanism outputs $\mu \in \Delta(V)$ at the end of period $t$, then $\mu$ is indeed the belief the principal holds about the agent at the end of that period. Motivated by Theorem 3.1, Section 3.2 studies the PBE of the *canonical game*. It follows immediately from Theorem 3.1 that any equilibrium payoff of the mechanism-selection game is also an equilibrium payoff of the canonical game, after adapting the strategy profiles and systems of beliefs to the canonical game. Because the canonical game has a smaller set of deviations for the principal than the mechanism-selection game, one may conjecture that there are dynamic mechanisms consistent with equilibrium in the canonical game, which would not be consistent with equilibrium in the mechanism-selection game. Proposition 3.1 shows this conjecture is false.

### 3.1  Revelation Principle for Sequentially Optimal Mechanism Design

**Theorem 3.1.** Fix any PBE of the mechanism-selection game, $\langle \Gamma^*, (\pi_v^*, r_v^*)_{v \in V}, \mu^* \rangle$. Then there exists a payoff-equivalent PBE, $\langle \Gamma', (\pi_v', r_v')_{v \in V}, \mu' \rangle$, such that

1. At all histories, the principal offers canonical mechanisms, that is, $(\forall h^t)(\forall \mathbf{M}_t : \Gamma'(h^t)(\mathbf{M}_t) > 0)$, $\mathbf{M}_t \in \mathcal{C}$.

2. At all histories where the principal assigns positive probability to the agent's type being $v$, the agent participates with probability 1 when her type is $v$, that is, $(\forall v, \forall h_A^t \in H_A^t)(\forall \mathbf{M}_t \in \text{supp } \Gamma'(h^t))$ $\pi_v'(h_A^t, \mathbf{M}_t) = 1$ whenever $\mu'(h^t)(v) > 0$.

3. At all histories, the agent reports her type truthfully, that is, $(\forall v, \forall h_A^t \in H_A^t)(\forall \mathbf{M}_t \in \text{supp } \Gamma'(h^t))$, $r_v'(h_A^t, \mathbf{M}_t, 1) = \delta_v$.

21

4. At all histories, recommended beliefs coincide with realized beliefs $t + 1$:

$$\mu'(h^t, \mathbf{M}_t, 1, \mu)(v) = \frac{\mu'(h^t)(v)\beta^{\mathbf{M}_t}(\mu|v)}{\sum_{v' \in V} \mu'(h^t)(v')\beta^{\mathbf{M}_t}(\mu|v')} = \mu(v).$$

The proof is in Appendix B. In what follows, we provide a sketch of the main steps in the proof.

The first main step shows that, without loss of generality, the agent's participation and reporting strategy conditions only on her type $v$ and the public history. This step, which follows from Proposition A.1 in Appendix A, is key to showing that the set of canonical input messages is the set of type reports. If the agent conditioned her strategy on the payoff-irrelevant part of her private history, the principal would need to elicit $h_A^t$ together with $v$ in order to replicate the agent's behavior in the mechanism.

We now qualify what we mean by *without loss of generality*: We show that given a PBE in which the agent conditions her strategy on the payoff-irrelevant part of her private history at some public history $h^t$, there exists another payoff-equivalent PBE in which she does not and in which the principal obtains the same payoff after each continuation history consistent with $h^t$ and the equilibrium strategy. The proof of this consists of two parts. First, we observe that because the input messages are payoff irrelevant and unobserved by the principal, if the agent chooses different strategies at $(v, h_A^t)$ and $(v, h_A^{t\,\prime})$, with $h_A^t, h_A^{t\,\prime} \in H_A^t(h^t)$, then she is indifferent between these two strategies. However, the principal may not be indifferent between these two strategies. The second part shows we can build an alternative strategy that does not condition on $h_A^t$ beyond $h^t$ and gives the principal the same continuation payoff.

This first step also gives us an important conceptual insight: The principal cannot peak into his past correlating devices. To do so, he would like to ask the agent to report to him what she did in the previous mechanisms. Proposition A.1 shows that this information cannot be elicited in any payoff-relevant way.

The second main step shows that, without loss of generality, there is a one-to-one mapping between the output messages used at history $h^t$ and the posterior beliefs

of the principal in the PBE at history $h^t$ (see Proposition A.3 in Appendix A). This step follows mainly from the observations we made in the discussion in Section 2.2. The principal may have two other uses for the output messages. On the one hand, because the allocation must be measurable with respect to the output messages, he may use them to offer a richer set of alternatives. On the other hand, he may use the output messages to coordinate continuation play. Proposition A.3 shows that randomized allocations and the access to the public correlating device can achieve these two goals, respectively.

These two steps deliver that, without loss of generality, input messages can be taken to be type reports and output messages can be taken to be the designer's beliefs about the agent's type. After all, knowing the agent's type is all that is needed to replicate her behavior within the mechanism, and hence the relevant beliefs for the principal are about the agent's payoff-relevant type.

Proposition A.2 in Appendix A shows that having the agent participate in the mechanism is without loss of generality (we discuss at the end of the section why Theorem 3.1 only requires this for types with positive probability.). The logic is similar to the one in the case of commitment: Whatever the agent obtains when she does not participate can be replicated by making her participate. However, there is a caveat: When the agent does not participate, her outcome is an allocation for today and a continuation mechanism for tomorrow. Therefore, we must guarantee that, when the agent participates, the principal still offers the same continuation as when she did not participate.

With these preliminary steps at hand, the proof of Theorem 3.1 in Appendix B shows that any mechanism $\mathbf{M}_t$ offered by the principal at history $h^t$ can be replicated by a canonical mechanism $\mathbf{M}_t^C = \langle (V, \beta^{\mathbf{M}_t^C}, \Delta(V)), \alpha^{\mathbf{M}_t^C} \rangle$ as follows. The second step implies that there is an invertible mapping which maps each output message into the belief over types that it induces:

$$\sigma(\mathbf{M}_t)(s_t) = \sum_{h_A^t \in H_A^t(h^t), m_t \in M^{\mathbf{M}_t}} \mu^*(h^t, \mathbf{M}_t, 1, s_t)(\cdot, h_A^t, \mathbf{M}_t, 1, m_t, s_t).$$

Note that we obtain the belief over $V$ by taking the marginal over all agent histories

consistent with the public history $h^t$. Using this, we can define a communication device $\beta^{\mathbf{M}_t^C} : V \mapsto \Delta^*(\Delta(V))$ and an allocation rule $\alpha^{\mathbf{M}_t^C} : \Delta(V) \mapsto \Delta^*(A)$ as follows:

$$\beta^{\mathbf{M}_t^C}(\mu|v) = \sum_{m \in M^{\mathbf{M}_t}} \beta^{\mathbf{M}_t}(\sigma^{-1}(\mathbf{M}_t)(\mu)|m) r_v^*(h_A^t, \mathbf{M}_t, 1)(m)$$

$$\alpha^{\mathbf{M}_t^C}(\mu) = \alpha^{\mathbf{M}_t}(\sigma^{-1}(\mathbf{M}_t)(\mu)).$$

The proof then shows that when faced with this mechanism, the agent's best response is to participate and report truthfully and that when the principal observes an output of $\mu$, his beliefs are indeed $\mu$.

We have yet to discuss why Theorem 3.1 only requires that the agent participates with probability 1 is required for her types to which the principal assigns positive probability. Consider then a history $h^t$ such that the principal's belief assigns probability 0 to the agent's type being $v^\star$. Suppose the principal selects mechanism $\mathbf{M}_t$. Assume also the agent's strategy at $v^\star$ specifies sending an input message $m^\star$, which assigns positive probability to an output $s^\star$. Finally, assume that $s^\star$ has zero probability under all other $m \in M^{\mathbf{M}_t}$. PBE does not impose restrictions on the principal's belief when he observes $s^\star$; in particular, it could be that $\mu^*(h^t, \mathbf{M}_t, 1, s^\star) = \mu^*(h^t, \mathbf{M}_t, 1, s')$, where $s'$ is an output message with positive probability under the equilibrium strategy. This would, of course, break the one-to-one mapping between output messages and posterior beliefs.

To deal with the aforementioned issue, we show that, given a PBE, we can always modify the mechanisms chosen in equilibrium by the principal so that the agent does not have access to messages like $m^\star$. Namely, he can make the distribution of the communication device for any such message the same as that of a message that is used on the path. The principal can always do this without affecting the incentives of those types that have positive probability; however, he may change the participation incentives of those types that have probability 0.

24

## 3.2 The Canonical Game

Theorem 3.1 shows that any equilibrium payoff of the game between the principal and the agent can be achieved with the principal selecting at each history a canonical mechanism such that the agent participates with probability one and reports her type truthfully. This observation motivates the analysis in this section where we study the equilibria of the canonical game.

An immediate corollary of Theorem 3.1 is the following:

**Corollary 3.1.** *Any PBE payoff of the mechanism-selection game can be achieved as a PBE payoff of the canonical game.*

Because it features a restricted set of choices for the principal, one may suspect that in the canonical game, the principal is able to implement more mechanisms than in the mechanism-selection game. However, this is not the case. Indeed, we show that given any equilibrium of the mechanism-selection game, without loss of generality, the *best* deviation for the principal after *any* history can be achieved by offering a canonical mechanism for that period and also in the continuation histories, whereas the agent participates with probability one and truthfully reports her type. This observation implies the canonical game contains all relevant deviations for the principal. It is not then possible to achieve payoffs in the canonical game that cannot be achieved in the mechanism-selection game. This is recorded in Proposition 3.1 below:

**Proposition 3.1.** *If $\langle \Gamma^{*^C}, (\pi_v^{*^C}, r_v^{*^C})_{v \in V}, \mu^{*^C} \rangle$ is a PBE of the canonical game, then there is an equilibrium of the mechanism-selection game $\langle \Gamma^*, (\pi_v^*, r_v^*)_{v \in V}, \mu^* \rangle$ that achieves the same payoff.*

Two important lessons follow from Proposition 3.1 and its proof. First, to characterize the equilibrium payoffs of the mechanism-selection game, it suffices to characterize the equilibrium payoffs of the canonical game.[16] Second, the proof of

---

[16]It is not obvious that such a result should hold. To see this, we reason by analogy with the informed principal problem of Myerson (1983), where the principal is also a player, and focusing on deviations to direct and incentive compatible mechanisms is with loss of generality. Two of the equilibrium notions Myerson analyzes have analogues in our paper. In *expectational equilibria*, the principal can choose from *any* mechanism, as in the model in Section 2. *Undominated mechanisms* are direct incentive compatible mechanisms that weakly dominate any other direct

Proposition 3.1 highlights that it is enough for the principal to look among those canonical mechanisms that incentivize the agent to participate and truthfully report her type. This second observation is important. In finite-horizon settings, it justifies writing down the principal's problem as a series of maximization problems subject to constraints: the participation and incentive compatibility constraints for the agent and the sequential rationality constraints for the principal. This provides a game-theoretic foundation for the programs studied by Bester and Strausz (2001) and Bester and Strausz (2007). In Doval and Skreta (2018b) we show how this approach also simplifies looking for the best equilibrium for the principal in an infinite-horizon problem.

## 4   Transferable utility and increasing differences

Section 4 considers a simplified version of the game in Section 2. The purpose is to show how one can harness the results in Section 3 to solve for the principal's optimal mechanism under limited commitment. In particular, our formulation of the canonical set of output messages as beliefs allows us to write the principal's problem as a *constrained* information design problem. Using this formulation and an extension[17] of the techniques in Le Treust and Tomala (2017), we characterize upper bounds on the set of posteriors used in an optimal mechanism. Along the way, we also highlight the differences between the problem considered here and the one introduced by Kamenica and Gentzkow (2011).

Consider the following simpler version of the game in Section 2:

- The agent observes her type $v_i \in \{v_1, \ldots, v_N\}$. Let $\mu_i^0 = Pr(v = v_i)$.

- The principal offers the agent a mechanism $\mathbf{M} = \langle (M, \beta, S), \alpha \rangle$.

- The agent observes $\mathbf{M}$ and decides whether to participate.

  – If she does not participate, $a^*$ is implemented.

  – If she participates, she privately submits a report $m \in M$:

---

and incentive compatible mechanism, as in the canonical game in Section 3.2. Myerson shows that strong solutions, which is a strengthening of undominated mechanisms, are expectational equilibria, but the reverse does not necessarily hold.

[17]See Doval and Skreta (2018a) for further details.

- $s \in S$ is drawn according to $\beta(\cdot|m)$, which is publicly observed,

- $a \in A$ is drawn according to $\alpha(\cdot|s)$, which is publicly observed.

- The principal selects an action $y \in Y(a)$, where $Y(a) \subseteq Y$ is a compact (possibly finite) space.

The *non-contractible* action $y$ captures in reduced form the principal's limited commitment: In the example in Section 2.1, $y$ corresponds to the choice of mechanism in the second period.[18] The correspondence $Y(a)$ plays the role of the correspondence $\mathcal{A}$ in Section 2: It captures how past allocations may affect the principal's available choices in the continuation.

The above is the game that underlies the maximization problem analyzed by Bester and Strausz (2007). Theorem 3.1 and Proposition 3.1 provide a game-theoretic foundation for why the search for the principal's best equilibrium can be cast in terms of such a program. To facilitate the comparison between the papers, we follow their notation as much as possible. In what follows, $w_i(a, y)$ denotes the principal's utility when $v = v_i$; similarly, $u_i(a, y)$ is the agent's utility when her type is $v_i$.

In standard mechanism design fashion, we focus here on the case of transferable utility and *increasing differences*, leaving the full analysis to Section I in the Supplementary Material.[19] First, we assume $A = Q \times \mathbb{R}$, where $q \in Q$ denotes the physical part of the allocation and $t \in \mathbb{R}$ denotes a monetary transfer from the agent to the principal. Hereafter, we take $Y(q, t) = Y(q)$, and in a slight abuse of notation, we denote $u_i(a, y) = u_i(q, y) - t, w_i(a, y) = w_i(q, y) + t$.[20] Second, because mechanisms in our setting determine lotteries over outcomes, the appropriate notion of increasing differences is the one in Kartik et al. (2017):

---

[18]The beauty of the simple model is that the non-contractible part of the allocation $y$ may stand for other forms of contractual incompleteness, such as renegotiation. From this point of view, we believe the techniques presented herein could be used to understand optimal contracting in other environments of interest, where distortions arise, for instance, from the need to renegotiate contracts or hold-up problems.

[19]We comment at the end of this section how the results we obtain translate to the case of non-transferable utility

[20]Under transferable utility, if the action $y$ represents the choice of a continuation mechanism, then the assumption that $Y(\cdot)$ does not depend on $t$ is innocuous. Section I in the supplementary material does not restrict how the correspondence $Y(\cdot)$ depends on $a$.

**Definition 4.1** (Kartik et al. (2017))**.** The family $\{u_i\}_{i=1}^N$ satisfies *monotonic expectational differences* if for any two distributions $P, Q \in \Delta(A \times Y)$ $\int u_i(\cdot)dP - \int u_i(\cdot)dQ$ is monotonic in $i$.[21]

The analysis in Section 3 implies the solutions to the following program characterize the PBE of the aforementioned game:[22]

$$\max_{\beta,\alpha,y} \sum_{i,h} \mu_i^0 \beta_{i,h}[\sum_{(q,t)} \alpha_h(q,t)(w_i(q, y_h(q)) + t)] \qquad (\mathcal{P})$$

$$s.t. \begin{cases} \sum_h \beta_{i,h} \sum_{(q,t)} \alpha_h(q,t)[u_i(q, y_h(q)) - t] \geqslant 0, i \geqslant 1 \\ \sum_h (\beta_{i,h} - \beta_{k,h}) \sum_{(q,t)} \alpha_h(q,t)[u_i(q, y_h(q)) - t] \geqslant 0, i, k \in \{1, 2, ..., N\}, i \neq k \\ y_h(q) \in y^*(\mu_h, q) \equiv \arg\max_{y \in Y(q)} \sum_{i=1}^N \mu_{h,i} w_i(q, y) \\ \mu_{h,i} \sum_{j=1}^N \mu_j^0 \beta_{j,h} = \mu_i^0 \beta_{i,h}, \end{cases}$$

where $\beta_{i,h} = \beta(\mu_h|v_i), \alpha_h = \alpha(\cdot|\mu_h)$ and $\mathcal{H} = \{1, \ldots, H\}$ indexes the posteriors. That is, the principal selects the best canonical mechanism from among the ones that (i) induce participation with probability 1, (ii) induce truthtelling with probability 1, and (iii) satisfy the principal's sequential rationality constraints. Implicit in this program is that the number of posteriors induced by the principal is also a variable of choice.

The rest of the section proceeds as follows. First, Proposition 4.1 shows how, under our assumptions, we can simplify the number of constraints in program $(\mathcal{P})$. Second, we show how to cast the simpler program as a constrained information design one. Finally, we use this connection to characterize an upper bound on the

---

[21]Kartik et al. (2017) show that $u$ satisfies monotonic expectational differences if, and only if, it takes the form $u_i(a, y) = g_1(a, y)f_1(i) + g_2(a, y) + c(i)$, where $g_1, g_2$ are finitely integrable and $f_1$ is monotonic.

[22]In the event that the agent does not participate in the mechanism, allocation $a^*$ gets implemented. Moreover, the principal chooses $y \in Y(a^*)$ to maximize his expected utility, where expectations are taken with respect to his beliefs after observing the agent does not participate of the mechanism, which is an off-path event. Implicit in the agent's participation constraint in program $(\mathcal{P})$ is that the above choices can be made so that the agent obtains a payoff of 0 *regardless* of her type. This allows us to focus on the issues of limited commitment, without having to worry of the issue of type dependent participation constraints.

number of posteriors in the optimal mechanism.

Transferable utility implies that focusing on mechanisms that do not randomize on transfers is without loss of generality. Hereafter, we replace $t$ with its expectation, denoted by $t_h$. Like increasing differences in mechanism design with commitment, monotonic expectational differences implies the solutions to ($\mathcal{P}$) coincide with the solutions to a simpler program, which imposes only a subset of the incentive compatibility constraints. Finally, both assumptions imply that the participation constraint of the lowest type binds. The above remarks are recorded in Proposition 4.1.

**Proposition 4.1.** *If* $\{u_i\}_{i=1}^N$ *satisfies monotonic expectational differences, then to characterize the solution to* ($\mathcal{P}$)*, it suffices to guarantee the following hold:*

1. *The agent's participation constraint binds when her type is* $v_1$.

2. *Adjacent incentive constraints are satisfied.*

See Appendix D for a proof. In mechanism design with commitment, we could simplify ($\mathcal{P}$) further by showing the *downward-looking*[23] incentive constraints always bind at the optimum. This then justifies the study of the so-called *relaxed program*:

$$\max_{\beta,\alpha,y} \sum_{i,h} \mu_i^0 \beta_{i,h} \Big[ \sum_q \alpha_h(q) w_i(q, y_h(q)) + t_h \Big] \qquad (\mathcal{R})$$

$$s.t. \begin{cases} \sum_h \beta_{1,h} [\sum \alpha_h(q) u_1(q, y_h(q)) - t_h] = 0, \\ \sum_h (\beta_{i,h} - \beta_{i-1,h})[\sum \alpha_h(q) u_i(q, y_h(q)) - t_h] = 0, i \in \{2, ..., N\} \\ y_h(q) \in y^*(\mu_h, q) \equiv \arg\max_{y \in Y(q)} \sum \mu_{h,i} w_i(q, y) \\ \Big[ \sum_{j=1}^N \mu_j^0 \beta_{j,h} \Big] \mu_{h,i} = \mu_i^0 \beta_{i,h} \end{cases},$$

---

[23]That is, the constraints that say $v_i$ does not report her type is $v_{i-1}$.

obtained by dropping the *monotonicity constraints*:[24]

$$\sum_h (\beta_{i,h} - \beta_{i-1,h}) \sum_q \alpha_h(q)(u_i(q, y_h(q)) - u_{i-1}(q, y_h(q))) \geqslant 0, i \in \{2, \ldots, N\}. \quad \text{(M)}$$

In mechanism design with commitment, it suffices to check that the solution to the relaxed program satisfies the monotonicity constraints, (M), to show it is the solution to $(\mathcal{P})$ (see the discussion in footnote 24).

However, in mechanism design with limited commitment, the solution to the relaxed program is not necessarily a solution to $(\mathcal{P})$ even if it satisfies the monotonicity constraints, when the type space is finite and there are three or more types. This is illustrated in Example 2 in Appendix D. Whereas in the relaxed program the binding downward-looking incentive constraints together with $v_1$'s participation constraint impose $N$ restrictions on the transfers $(t_h)$, the solution to the relaxed program might use less than $N$ posteriors. Therefore, finding transfers $t_h$ that satisfy all constraints may not possible.[25] Alternatively, not all downward-looking constraints may bind in the optimal mechanism.

Fortunately, the above is not an issue when there are two types or a continuum of types. In both cases, it is possible to show that downward looking constraints bind (see Appendices D and IV). Because most of the literature focuses on one of these cases, and because the relaxed program provides a useful benchmark, the rest of this section studies its properties.

We can use the binding constraints to substitute the transfers out of the principal's program and obtain the following:

---

[24]The constraints in equation (M) are obtained from combining the restriction that $v_i$ does not want to report $v_{i-1}$ and $v_{i-1}$ does not want to report $v_i$. Under Definition 4.1, the binding downward-looking incentive constraints together with the monotonicity constraints imply the local constraints in Proposition 4.1.

[25]This is never an issue in mechanism design with commitment: Without loss of generality, we can always have one transfer for each type.

$$\max_{\tau,\alpha,y} \sum_h \tau(\mu_h) \sum_{i=1}^N \mu_{h,i} \sum_q \alpha_h(q)\left[w_i(q, y_h(q)) + u_{i,h} - \frac{1 - \sum_{n\leqslant i} \mu_n^0}{\mu_i^0}(u_{i+1,h} - u_{i,h})\right]$$
$$\text{s.t.} \sum \tau(\mu_h)\mu_h = \mu^0,$$

where $\tau(\mu_h) = \sum_i \mu_i^0 \beta_{i,h}$ and $u_{i,h} = u_i(q, y_h(q))$.

Define:

$$\hat{u}_i(q, y; \mu^0) \equiv u_i(q, y) - \frac{1 - \sum_{n\leqslant i} \mu_n^0}{\mu_i^0}(u_{i+1}(q, y) - u_i(q, y)),$$

$$w_i(\alpha, y_\mu(\alpha)) + \hat{u}_i(\alpha, y_\mu(\alpha); \mu^0) \equiv \sum_q \alpha(q)(w_i(q, y_\mu(q)) + \hat{u}_i(q, y_\mu(q); \mu^0))$$

$$\tilde{w}(\alpha, \mu; \mu^0) \equiv \mathbb{E}_\mu\left[w_i(\alpha, y_\mu(\alpha))) + \hat{u}_i(\alpha, y_\mu(\alpha); \mu^0)\right],$$

where $\hat{u}_i$ is type $i$'s virtual utility from $(q, y)$ and $\tilde{w}$ is the expectation according to $\mu$ of the virtual surplus at $\{\alpha, y_\mu(\cdot))\}$, for some selection $y_\mu(q) \in y^*(\mu, q)$ (see Remark 4.1). Moreover, we drop the index $h$ because thinking about these objects as functions of beliefs $\mu$ in what follows is useful.

Program ($\mathcal{R}$) is then equivalent to

$$\max_{\tau,\alpha,y} \mathbb{E}_\tau \tilde{w}(\alpha, \mu; \mu^0) \tag{3}$$
$$\text{s.t. } \mathbb{E}_\tau \mu = \mu^0$$

That is, the solution to the relaxed problem is obtained by maximizing a version of the virtual surplus, represented by $\tilde{w}$, and then choosing a distribution over posteriors that averages out to the prior. Equation (3) generalizes the program obtained in Section 2.1. The following remark is in order:

**Remark 4.1** (Tie-breaking in favor of the principal)**.** So far, we have remained silent about how $y_\mu(q)$ is chosen, beyond the restriction that $y_\mu(q) \in y^*(\mu, q)$. We can use the function $\tilde{w}(q, \mu; \mu^0)$ to determine how to break the possible ties in $y^*$ and make the principal's objective upper-hemicontinuous. In fact, if $y, y' \in$

$y^*(\mu, q)$, then in the relaxed problem, $y$ is selected as long as

$$\mathbb{E}_\mu[w_i(q, y) + \hat{u}_i(q, y; \mu^0)] \geqslant \mathbb{E}_\mu[w_i(q, y') + \hat{u}_i(q, y'; \mu^0)].$$

In other words, ties are broken in favor of the virtual surplus.

We now illustrate how to solve the program in (3). Towards this, fix the selection $y$ as in Remark 4.1. Because the program is separable in the allocation, $\alpha$, across posteriors, the solution can be obtained in two steps. First, for each posterior $\mu$, we maximize $\tilde{w}(\cdot, \mu; \mu^0)$ with respect to $\alpha$. Denote the value of this problem $\hat{w}(\mu; \mu_0)$. Second, we choose $\tau$ to maximize $\hat{w}(\mu; \mu_0)$ subject to the constraint that the posteriors must average out to the prior, $\mu^0$. This separability between the choice of the allocation rule, $\alpha$, and the communication device, $\beta$, is afforded by ignoring the monotonicity constraints in (M). The latter may impose additional restrictions on how the allocation varies across different posteriors.

This discussion implies the solution to (3) can be obtained by solving:

$$\max_\tau \mathbb{E}_\tau \underbrace{\max_\alpha \tilde{w}(\alpha, \mu; \mu^0)}_{\hat{w}(\mu; \mu^0)} \tag{4}$$
$$\text{s.t. } \mathbb{E}_\tau \mu = \mu^0$$

An advantage of the above formulation is that a straightforward application of Carathéodory's theorem (see Rockafellar (1970)) implies that in (4), the solution never uses more than $N$ posteriors:[26]

**Proposition 4.2.** *The solution to ($\mathcal{R}$) uses at most $N$ posteriors.*

Then, if the solution to the relaxed program satisfies the monotonicity constraints *and* it is possible to find transfers $(t_h)$ that satisfy the downward looking binding incentive constraints, we have found a solution to the principal's problem, ($\mathcal{P}$).

In many instances, however, the solution to ($\mathcal{R}$) will fail to satisfy the mono-

---

[26]Bester and Strausz (2007) derive this result using methods in semi-infinite linear programming.

tonicity constraints, (M). As we show next, adding as many posteriors as binding monotonicity constraints at the optimum may be necessary:

**Proposition 4.3.** *Consider the program obtained by adding the monotonicity constraints* (M) *to the relaxed program* ($\mathcal{R}$). *The solution to the new program uses at most $N + K$ posteriors, where $K$ is the number of binding constraints at the optimum.*

The proof is in Appendix D and follows from extending the techniques in Le Treust and Tomala (2017) to our setting, where we have multiple inequality constraints and equality constraints.

Finally, we note the connection between our problem and a constrained information design problem holds beyond the case of transferable utility, as illustrated in Section I in the supplementary material. In particular, we show the assumption of monotonic expectational differences also reduces the problem to the analysis of the local incentive constraints. Moreover, we can again bound the number of posteriors by $3N - 1$.

Whereas the above formulation harnesses the connection between our problem and the one studied in information design, we close the section by highlighting two conceptual differences with this literature. The reader eager to see the results in the next section can skip it without loss of continuity.

First, the function $\hat{w}(\mu; \mu^0)$ in equation (4) stands for the sender's objective function, $\hat{v}(\mu)$, in Kamenica and Gentzkow (2011). Recall that in Kamenica and Gentzkow (2011), $\hat{v}(\mu)$ is the sender's expected utility of the receiver's optimal action when the posterior is $\mu$, where expectations are taken with respect to $\mu$. Two differences are worth pointing out. First, in our setting, the first-period principal (the sender in Kamenica and Gentzkow (2011)) also takes an action for each posterior he induces, because he chooses the allocation $\alpha$.

Second, the principal's objective function in equation (3) depends both on the posterior, $\mu$, and the prior, $\mu_0$, whereas in Kamenica and Gentzkow (2011), it only depends on the posterior. We already saw an instance of this in the example studied in Section 2.1. In fact, we saw that the virtual values are calculated using

the prior distribution, because this distribution is the one the principal uses to calculate the probability with which he leaves rents to the different types of the agent. That the principal's payoff depends both on the prior *and* the posterior may come as a surprise because when a distribution $F$ can be written as a convex combination of distributions $F_s$, so that $F = \sum_{s=1}^{S} \lambda_s F_s$, then

$$\int \frac{(1-F)}{f} dF = \sum_{s=1}^{S} \lambda_s \int \frac{1-F_s}{f_s} dF_s.^{27}$$

That is, the posterior *information handicap* averages out to the prior information handicap. Thus, we may have expected that the information handicap in $\tilde{w}(\alpha, \mu; \mu_0)$ could be written solely as a function of $\mu$. Only when the allocation is the same for all induced posteriors, and hence no type of the agent obtains rents, we can think of the principal's objective as only depending on the posterior.

## 5 Discussion

### 5.1 Recommendations as output messages

As discussed in the introduction, in the finite-horizon case, there is another candidate for a canonical language: In period $t$, each output message could be associated to an allocation for period $t$ and a recommended allocation for the subsequent periods. We use the formulation in Section 4 to discuss this formally.

Section 4 illustrated how the relaxed program can be formulated as an information design problem, where the principal in period 1 designs both the allocation for the agent and the information structure for the principal in period 2 (see equation (3)).

We borrow the terminology in Kamenica and Gentzkow (2011) and say a mechanism is *straightforward* if $S \subseteq \cup_{q \in Q} \{q\} \times Y(q)$ and after message $s = (q_s, y_s)$, the principal chooses $q_s$ in period 1 to maximize $\tilde{w}(\alpha, \mu_s; \mu^0)$ and $y_s$ in period 2 to maximize $\sum \mu_{i,s} w_i(q_s, y)$ in period 2, where $\mu_s$ are the beliefs generated by output message $s$.

---

**Proposition 5.1.** *The following are equivalent:*

1. *There exists a mechanism $\langle (V, \beta, S), \alpha \rangle$ and a mapping $y : S \times Q \mapsto \cup_{q \in Q} Y(q)$ that solves $(\mathcal{R})$.*

2. *There exists a mechanism $\langle (V, \beta, \Delta(V)), \alpha \rangle$ and a mapping $y : \Delta(V) \times Q \mapsto \cup_{q \in Q} Y(q)$ that solves $(\mathcal{R})$.*

3. *There exists a straightforward mechanism that solves $(\mathcal{R})$.*

The proof is in Appendix E. Item 3 highlights that in the relaxed program $(\mathcal{R})$, the set of output messages can also be taken to be recommendations for both incarnations of the principal, as in sender-receiver models of information design. The proposition uses both the separability between the allocation, $\alpha$, and the information policy, $\tau$, discussed in Section 4 and that there is a final period in which the principal takes the non-committed action, $y$.[28] The separability guarantees the first-period principal chooses to implement allocations $q$ that are optimal for each posterior he induces for period 2. However, in the case of infinite horizon, the language of recommendations is *self-referential*: The principal would need to recommend the continuation mechanisms, which themselves involve a set of output messages. Thus, an advantage of the approach we advocate is that we can always resort to beliefs regardless of the game at hand.

## 5.2 Implementation via contracts

Section 5.2 characterizes within the environment of Section 4 the tuples $(\beta, q, y)$ that can be implemented using the contracts studied previously in the literature. An important difference between the mechanisms used by Hart and Tirole (1988), Laffont and Tirole (1988), Freixas et al. (1985), and Bester and Strausz (2001) and the ones considered here is that whereas in the former papers, the principal observes the agent's choice out of a menu, in the latter, the agent's input into the communication device is not observed. A consequence of this is that in the former setting, the agent has to be indifferent between all the elements of the menu

---

[28]If the final period corresponded to a design problem, such as the sale of a durable good example, one could resort to the revelation principle in Myerson (1982) to reduce the principal's actions in the final period to the induced allocations.

that she selects with positive probability. By contrast, in our setting, the agent's incentive compatibility constraint has to hold in expectation: Although she may not be indifferent between all the allocations that arise with positive probability after she communicates with the mechanism, on average, they must be better than what she would obtain by reporting any other type.

Fix a *canonical* communication device $\beta : V \mapsto \Delta^*(\Delta(V))$ and a tuple $(q, y) : \Delta(V) \mapsto Q \times Y$, where we denote by $y(\mu) = y(q(\mu), \mu)$.[29] We want to find $t' : \Delta(V) \mapsto \mathbb{R}$ such that for all $v_i \in V$, for all $\mu : \beta(\mu|v_i) > 0$, and for all $\mu' : \sum_{v'} \beta(\mu'|v') > 0$,

$$u_i(q(\mu), y(\mu)) - t'(\mu) \geqslant u_i(q(\mu'), y(\mu')) - t'(\mu'). \qquad \text{(DIC-P)}$$

Note equation (DIC-P) corresponds to the definition of equivalence in Mookherjee and Reichelstein (1992).[30] Indeed, the problem we intend to solve is similar in spirit to the one analyzed in the literature that studies the equivalence between Bayesian and dominant-strategy implementation (Manelli and Vincent (2010) and Gershkov et al. (2013)). However, there are some differences that, although subtle, turn out to have important implications. First, in that literature, this problem only makes sense when there are multiple agents, whereas in our case, the source of randomness the agent faces (the randomization by the communication device) is endogenously chosen by the principal. Second, allocation and transfers in that setting depend on the agent's type, whereas here they depend on the belief induced when the principal observes the output message. As we show next, this implies joint restrictions on the communication device and the allocation rule.

As in Section 4, we assume the agent's preferences satisfy monotonic expectational differences. Thus, label the types so that $v_1 < \cdots < v_N$. Note that if $u_i(q, y)$ satisfies Definition 4.1, then $u_i(q, y)$ has increasing differences. In effect, $u_i(q', y') - u_i(q, y) = f(i)(g_1(q', y') - g_1(q, y)) + g_2(q', y') - g_2(q, y)$, which is strictly

---

[29]To keep notation simple, we ignore the possibility that $q$ and $y$ may be randomized allocations. It is immediate that this restriction is not necessary for the results.

[30]Mookherjee and Reichelstein (1992) also require that $t'$ raises the same revenue as $t$ (see equation (DIC-T) below). We focus for now on the possibility of guaranteeing (DIC-P) holds and discuss the difficulties associated with guaranteeing the same revenue is collected at the end of Section 5.2.

increasing in $i$ as long as $g_1(q', y') - g_1(q, y) \neq 0$. In what follows, we make the following assumption:

**Assumption 1.** For all $\mu, \mu' \in \Delta(V)$ such that $(\sum_{v \in V} \beta(\mu|v)) \times (\sum_{v \in V} \beta(\mu'|v)) > 0$, we assume $g_1(q(\mu), y(\mu)) - g_1(q(\mu'), y(\mu')) \neq 0$.

Given two beliefs $\mu$ and $\mu'$, let

$$D_i(\mu, \mu') \equiv u_i(q(\mu), y(\mu)) - u_i(q(\mu'), q(\mu'))$$

denote the difference in payoffs from $(q(\mu), y(\mu))$ and $(q(\mu'), y(\mu'))$, when the agent type is $v_i$. The content of Assumption 1 is that $D_i(\mu, \mu')$ is strictly increasing in $i$.

We have the following:

**Proposition 5.2.** *Suppose the agent's Bernoulli utility function satisfies Definition 4.1 and $(\beta, q, y)$ satisfies Assumption 1. A necessary and sufficient condition for $(\beta, q, y)$ to satisfy (DIC-P) is that $(\beta, q, y)$ satisfies*

1. *For all $i \in \{1, \ldots, N\}$ and $j > i$, for all $\mu$ and $\mu'$,*

$$D_j(\mu', \mu) \geqslant D_i(\mu', \mu), \tag{DIC-M}$$

   *whenever $\mu'(v_j)\mu(v_i) > 0$.*

2. *$\beta$ induces a* monotone information structure*: We can label the beliefs induced by $(\beta, q, y)$, $\{\mu_1, \ldots, \mu_M\}$ so that*

   (a) *If $i < j$, then $\max supp \ \mu_i \leqslant \min supp \ \mu_j$,*

   (b) *For any $i$, there are at most three beliefs $\{\mu_i, \mu_{i+1}, \mu_{i+2}\}$ for which $v_i$ has positive probability. Moreover, if $v_i$ has positive probability in all three, then $\mu_{i+1}(v_i) = 1$.*

The proof is in Appendix F. The first condition is the equivalent to the standard monotonicity condition for dominant-strategy incentive compatibility: For any two beliefs $\mu$ and $\mu'$, the utility differential of the allocations $(q, y)$ induced at these beliefs is higher for higher types. The second is new to our setting. Recall that, under Assumption 1, $D_i(\mu, \mu')$ is strictly increasing in $i$, which places restrictions

on the support of the beliefs induced by the principal's mechanism. In particular, to satisfy (DIC-P), when the agent's type is $v_i$, she must be indifferent between the allocations $(q, y)$ that correspond to posteriors that assign positive probability to $v_i$. Monotonic expectational differences implies that when $v_i$ is indifferent between $(q(\mu), y(\mu))$ and $(q(\mu'), y(\mu'))$ and $D_i(\mu, \mu') \geqslant 0$, then all types higher than $v_i$ (weakly) prefer $(q(\mu), y(\mu))$ to $(q(\mu'), y(\mu'))$ (and the opposite holds for lower types). Assumption 1 then guarantees that higher types cannot be in the support of $\mu'$ (and the opposite holds for lower types).

Note that we cannot dispense with the assumption that $\beta$ induces a monotone information structure. In the example in Appendix B in Bester and Strausz (2007), the agent can be of one of two types and her utility satisfies Definition 4.1; however, the optimal mechanism induces three posteriors under which both types have positive probability and having both types be indifferent between the three allocations induced by the mechanism is not possible.

Besides allowing us to connect the results in this paper with the previous literature on mechanism design with limited commitment, the result in Proposition 5.2 is also of practical value. Section 4 highlights that the characterization of $S$ as the set of beliefs over the agent's type allows us to harness the tools of mechanism and information design to solve for the principal's optimal mechanism. Proposition 5.2 allows us then to check when the solution to the program in Section 4 is also a solution to the problem in which the principal observes the agent's choice out of a menu.

Proposition 5.2 is silent about whether the transfer scheme $t'$ collects the same revenue as the original mechanism did, that is, whether $t'$ also satisfies that

$$\sum_{v \in V} \beta(\mu|v)t'(\mu) = \sum_{v \in V} \beta(\mu|v)t(\mu), \qquad \text{(DIC-T)}$$

for all $\mu$ such that $\sum_{v \in V} \beta(\mu|v) > 0$.

Contrary to Mookherjee and Reichelstein (1992), we do not find that once (DIC-P) is satisfied, then (DIC-T) is satisfied. In particular, with two types, only when the solution features two beliefs, $\underline{\mu}(v_2) < \mu^0 < \overline{\mu}(v_2) = 1$, can one satisfy both (DIC-P)

and (DIC-T). We leave for future research the full analysis of the conditions under which both (DIC-P) and (DIC-T) hold.[31]

### 5.3 Multiple agents: Example in *Bester and Strausz (2000)*

Bester and Strausz (2000) show that, with multiple agents, the result for the single-agent case in Bester and Strausz (2001) no longer holds. That is, if $M = S$ and $\beta$ is deterministic, then there are equilibria with mechanisms in which $M \neq V$, whose payoffs cannot be replicated with canonical mechanisms. To keep the discussion self-contained, we replicate here their example and then explain why it does not invalidate the possibility of extending our techniques to the multi-agent case.

**Example 1** (Bester and Strausz (2000))**.** There are two agents, labeled 1 and 2. Only agent 1 has private information; let $v \in \{\underline{v}, \overline{v}\}$ denote her type. The prior that the type is $\underline{v}$ is denoted by $\mu \in [0, 1]$. The set of allocations $A = [0, 2]$. Payoffs are defined as follows:

$$W(a, \underline{v}) = -a^2, W(a, \overline{v}) = -(2 - a)^2$$
$$U_1(a, \underline{v}) = -(0.5 - a)^2, U_1(a, \overline{v}) = -(1.5 - a)^2$$
$$U_2(a) = -10(1 - a)^2.$$

That is, the principal's payoff depends on agent 1's type, whereas agent 2's payoff does not.

The timing is as follows. The principal selects a communication device for the agents, who then submit messages. Upon seeing the messages, the principal chooses $a \in A$.

The principal's payoffs are such that if, after seeing $m$, his posterior is $\mu(m)$, then he chooses allocation

$$a^*(m) = 2(1 - \mu(m)).$$

---

[31]Note that Mookherjee and Reichelstein (1992), Manelli and Vincent (2010), Gershkov et al. (2013) use the representation of the agent's utility function obtained via the envelope theorem to pin down transfers. As discussed in Section 4, we cannot guarantee downward-looking constraints bind at the optimum when there are three or more types.

Bester and Strausz (2000) construct an equilibrium with three messages $\{m_a, m_b, m_c\}$ that cannot be replicated with messages $\{\underline{v}, \overline{v}\}$. Let $M = \{m_a, m_b, m_c\}$. Then there is a PBE such that

$$\mu(m_a) = 1, \qquad \mu(m_b) = 1/2, \qquad \mu(m_c) = 0$$
$$a^*(m_a) = 0, \qquad a^*(m_b) = 1, \qquad a^*(m_c) = 2.$$

However, the above mechanism cannot be replicated by a mechanism with two messages, when the principal observes the output messages. The reason is not that agent 1 is not indifferent between the allocations he obtains at different messages, so that he is not willing to carry out the randomization himself. Rather, under the assumptions of Bester and Strausz (2001), the mechanism can only have as many input messages as output messages. Because the agent may be of one of two types, the mechanism can only have two input messages. Therefore, the agent does not have enough room to do the mixing and generate the required posteriors for the principal.

However, if we allow the principal to offer canonical mechanisms as the ones in this paper, the following communication device implements the same allocation as the non-canonical mechanism:

$$\beta(\mu(m_a)|\underline{v}) = 1/2, \beta(\mu(m_a)|\overline{v}) = 0,$$
$$\beta(\mu(m_b)|\underline{v}) = 1/2, \beta(\mu(m_b)|\overline{v}) = 1/2,$$
$$\beta(\mu(m_c)|\underline{v}) = 0, \beta(\mu(m_c)|\overline{v}) = 1/2.$$

After all, the result that taking $M \simeq V$ and $S \simeq \Delta(V)$ is without loss of generality dispenses with the restriction that the number of input messages must coincide with the number of output messages.

Extending the model in Section 2 to the case of multiple agents involves dealing with a number of subtleties that merit a full separate study and are thus beyond the scope of this paper. We plan to address this important extension in future research.

## A  Proof of preliminary results

Appendix A is organized as follows:

**Proposition A.1** shows we can focus without loss of generality on equilibria of the game in which the agent does not condition her strategy on the payoff-irrelevant part of her private history.

**Proposition A.2** shows we can focus without loss of generality on equilibria of the game in which the agent participates with probability one.

**Proposition A.3** shows we can focus without loss of generality on equilibria in which there is a one-to-one mapping between the output messages generated in the mechanism and the continuation beliefs the principal holds about the agent's type.

Because we have to deal with an abstract dynamic game, the proof is notationally involved. As a first pass to most results, except Proposition A.1, the reader is invited to first take a look at Appendix I in the supplementary material, where the constructions are performed in a two-period version of the model.

We need a few more pieces of notation and definitions.

First, as noted in footnote 12, some output messages can never be observed. Given a mechanism, $\mathbf{M}_t$, define $S^{*\mathbf{M_t}} = \{s \in S^{\mathbf{M}_t} : (\exists m \in M^{\mathbf{M}_t})\beta^{\mathbf{M}_t}(s|m) > 0\}$. Since removing public histories from the tree that are consistent with mechanism $\mathbf{M}_t$ and $s \in S^{\mathbf{M}_t} \backslash S^{*\mathbf{M}_t}$ does not change the set of equilibrium outcomes, hereafter, these histories are removed from the tree.

Second, fix a PBE of the dynamic mechanism-selection game $\langle \Gamma^*, (\pi_v^*, r_v^*)_{v \in V}, \mu^* \rangle$ and a public history $h^t$ for some $t \geqslant 0$ (if $t = 0$, then $h^0 = \varnothing$ denotes the initial public history). For $T \geqslant t$, the history $h^T = (h^t, \mathbf{M}_t, p_t, s_t, a_t, \ldots, \mathbf{M}_{T-1}, p_{T-1}, s_{T-1}, a_{T-1}, \omega_T)$ is on the path of the equilibrium strategy profile given $h^t$ if for all $t \leqslant \tau \leqslant T - 1$

$$h^{\tau+1} = (h^\tau, \mathbf{M}_\tau, p_\tau, s_\tau, a_\tau, \omega_{\tau+1}),$$

where

$$\mathbf{M}_\tau \in \text{supp } \Gamma^*(h^\tau)$$

$$\sum_{(v, h_A^\tau)} \mu^*(h^\tau)(v, h_A^\tau) \pi_v^*(h_A^\tau, \mathbf{M}_\tau)(p_\tau) > 0$$

$$p_\tau = 1 \Rightarrow \sum_{(v, h_A^\tau): \pi_v^*(h_A^\tau, \mathbf{M}_\tau) > 0} \mu^*(h^\tau)(v, h_A^\tau) \sum_{m \in M^{\mathbf{M}_\tau}} r_v^*(h_A^\tau, \mathbf{M}_\tau, 1)(m) \beta^{\mathbf{M}_\tau}(s_\tau | m) \alpha^{\mathbf{M}_\tau}(a_\tau | s_\tau) > 0$$

$$p_\tau = 0 \Rightarrow s_\tau = \varnothing, a_\tau = a^*.$$

That is, along the path from $h^t$ to $h^{\tau+1}$, the principal made choices according to his equilibrium strategy, the agent made participation choices according to her equilibrium strategy, and the output-message labels and allocations correspond to those in the mechanism chosen by the principal. Note that we do not say anything about the reports of the agent, because they are not part of the public history.

We sometimes need to talk about the histories that are on the path of the equilibrium strategy profile given a public history $h^t$ and a mechanism $\mathbf{M}_t$. The only difference with the above definition is that $\mathbf{M}_t$ need not have positive probability according to $\Gamma^*(h^t)$.

The above notation is used as follows. Proposition A.1 shows that for any PBE, there is a payoff-equivalent PBE in which the agent does not condition her strategy on the payoff-irrelevant part of her private history. To do so, starting from any history $h^t$, we need to modify the strategy *for all continuation histories on the path of the strategy*. Similarly, the main theorem shows we can transform any PBE of the game into one in which the principal's strategy selects only canonical mechanisms on and off the equilibrium path. To do so, we must map the continuation strategy starting from any history $h^t$ in the game to one in which the continuation strategy only offers canonical mechanisms. When we perform this mapping for $h^0$, we are doing the transformation for the path of the equilibrium strategy; when we do this for histories that can be reached from $h^t$, we are doing the transformation *for the path of the continuation strategy*.

**Proposition A.1.** *Fix a PBE $\langle \Gamma^*, (\pi_v^*, r_v^*)_{v \in V}, \mu^* \rangle$ and a public history $h^t$. Then, there exists a continuation strategy for the agent $(\pi_v^{**}, r_v^{**})_{v \in V}$ such that:*

1. *For any mechanism $\mathbf{M}_t$, for all $h_A^t, {h_A^t}' \in H_A^t(h^t)$, and for all $v \in V$, $\pi_v^{**}(h_A^t, \mathbf{M}_t) = \pi_v^{**}({h_A^t}', \mathbf{M}_t)$ and $r_v^{**}(h_A^t, \mathbf{M}_t, 1) = r_v^{**}({h_A^t}', \mathbf{M}_t, 1)$.*

2. *For all mechanisms $\mathbf{M}_t$, for $\tau \geqslant t+1$, for all histories $h^\tau$ on the equilibrium path starting from $(h^t, \mathbf{M}_t)$, for all $h_A^\tau, h_A^{\tau\,\prime} \in H_A^\tau(h^\tau)$, for all $\mathbf{M}_\tau \in supp\ \Gamma^*(h^\tau)$, and for all $v \in V$ $\pi_v^{**}(h_A^\tau, \mathbf{M}_\tau) = \pi_v^{**}(h_A^{\tau\,\prime}, \mathbf{M}_\tau)$ and $r_v^{**}(h_A^\tau, \mathbf{M}_\tau, 1) = r_v^{**}(h_A^{\tau\,\prime}, \mathbf{M}_\tau, 1)$.*

3. *For all histories $h^\tau$ on the equilibrium path starting from $h^t$, the continuation payoff for the principal at $(\Gamma^*, (\pi_v^{**}, r_v^{**})_{v \in V}, \mu^*)$ is the same as at $(\Gamma^*, (\pi_v^*, r_v^*)_{v \in V}, \mu^*)$; similarly, for each $v \in V$ and each $h_A^\tau \in H_A^\tau(h^\tau)$ the continuation payoff for the agent at $(v, h_A^\tau)$ is the same under both strategy profiles.*

4. $(\Gamma^*, (\pi_v^{**}, r_v^{**})_{v \in V}, \mu^*)$ *is also a PBE.*

*Proof.* Fix a PBE $\langle \Gamma^*, (\pi_v^*, r_v^*)_{v \in V}, \mu^* \rangle$. Let $h^t$ be a public history such that there exists $\mathbf{M}_t, h_A^t, h_A^{\prime t}$, both consistent with $h^t$, and $v \in V$ such that either $\pi_v^*(h_A^t, \mathbf{M}_t) \neq \pi_v^*(h_A^{t\,\prime}, \mathbf{M}_t)$ or $r_v^*(h_A^t, \mathbf{M}_t, 1) \neq r_v(h_A^{\prime t}, \mathbf{M}_t, 1)$.

Note that for each $m \in M^{\mathbf{M}_t}$, the agent's continuation payoff at $(v, h_A^t)$ and $(v, h_A^{t\,\prime})$ must be the same: after all, the continuation strategy of $(v, h_A^t)$ is feasible for $(v, h_A^{\prime t})$ and vice versa. Therefore, the agent at $(v, h_A^t)$ is not only indifferent between all the messages in the support of $r_v^*(h_A^t, \mathbf{M}_t, 1)$, but is also indifferent between all messages in the support of $r_v^*(h_A^{\prime t}, \mathbf{M}_t, 1)$. Therefore, the agent at $(v, h_A^t)$ is indifferent between $r_v^*(h_A^t, \mathbf{M}_t, 1)$ and any randomization between $r_v^*(h_A^t, \mathbf{M}_t, 1)$ and $r_v^*(h_A^{\prime t}, \mathbf{M}_t, 1)$.

Moreover, the above is true for any continuation public history that is reached with positive probability from $h^t$ for the same reasons. That is, for any $\tau \geqslant t$ and $h^\tau$ that succeeds $(h^t, \mathbf{M}_t)$ along which the principal follows $\Gamma^*$ and the agent at $(v, h_A^t)$ and at $(v, h_A^{\prime t})$ follows $\pi_v^*, r_v^*$ and for any $h_A^\tau, h_A^{\prime \tau}$ that succeed $h_A^t$ and $h_A^{\prime t}$, respectively, the agent is indifferent between her reporting strategy $r_v^*(h_A^\tau, \mathbf{M}_\tau, 1)$ and $r_v^*(h_A^{\prime \tau}, \mathbf{M}_\tau, 1)$ for $\mathbf{M}_\tau \in \Gamma^*(h^\tau)$.

Therefore, starting from $h^t$, the following is also an optimal strategy for the agent when her valuation is $v$. Consider first those types $v$ for which $\sum_{h_A^t \in H_A^t(h^t)} \mu^*(h^t)(v, h_A^t) > 0$. For any $h_A^t \in H_A^t(h^t)$, she participates with probability

$$\pi_v^{**}(h_A^t, \mathbf{M}_t) = \sum_{h_A^{t\,\prime}} \frac{\mu^*(h^t)(v, h_A^{t\,\prime})}{\sum_{\widetilde{h_A^t} \in H_A^t(h^t)} \mu^*(h^t)(v, \widetilde{h_A^t})} \pi_v^*(h_A^{t\,\prime}, \mathbf{M}_t),$$

and reports $m \in M^{\mathbf{M}_t}$ with probability:

$$r_v^{**}(h_A^t, \mathbf{M}_t, 1)(m) = \sum_{h_A^{t\prime} \in H_A^t(h^t)} \frac{\mu^*(h^t)(v, h_A^{t\prime}) \pi_v^*(h_A^{t\prime}, \mathbf{M}_t)}{\sum_{\widetilde{h_A^t} \in H_A^t(h^t)} \mu^*(h^t)(v, \widetilde{h_A^t}) \pi_v^*(\widetilde{h_A^t}, \mathbf{M}_t)} r_v^*(h_A^{t\prime}, \mathbf{M}_t, 1)(m),$$

as long as $\sum_{\widetilde{h_A^t} \in H_A^t(h^t)} \mu^*(h^t)(v, \widetilde{h_A^t}) \pi_v^*(\widetilde{h_A^t}, \mathbf{M}_t) > 0.$[32] For each $h^\tau$ that has positive probability from $(h^t, \mathbf{M}_t)$ and $h_A^\tau \in H_A^\tau(h^\tau)$, for each $\mathbf{M}_\tau \in \operatorname{supp} \Gamma^*(h^\tau)$, the agent participates with probability

$$\pi_v^{**}(h_A^\tau, \mathbf{M}_\tau) = \sum_{h_A^{\tau\prime} \in H_A^\tau(h^\tau)} \frac{\mu^*(h^\tau)(v, h_A^{\tau\prime})}{\sum_{\widetilde{h_A^\tau} \in H_A^\tau(h^\tau)} \mu^*(h^\tau)(v, \widetilde{h_A^\tau})} \pi_v^*(h_A^{\tau\prime}, \mathbf{M}_\tau), \tag{5}$$

as long as $\sum_{\widetilde{h_A^\tau} \in H_A^\tau(h^\tau)} \mu^*(h^\tau)(v, \widetilde{h_A^\tau}) > 0$ and reports $m \in M^{\mathbf{M}_\tau}$ with probability

$$r_v^{**}(h_A^\tau, \mathbf{M}_\tau, 1)(m) = \sum_{h_A^{\tau\prime} \in H_A^\tau(h^\tau)} \frac{\mu^*(h^\tau)(v, h_A^{\tau\prime}) \pi_v^*(h_A^{\tau\prime}, \mathbf{M}_\tau)}{\sum_{\bar{h}_A^\tau \in H_A^\tau(h^\tau)} \mu^*(h^\tau)(v, \bar{h}_A^\tau) \pi_v^*(\widetilde{h_A^\tau}, \mathbf{M}_\tau)} r_v^*(h_A^{\tau\prime}, \mathbf{M}_\tau, 1)(m),$$

$$\tag{6}$$

as long as $\sum_{\bar{h}_A^\tau \in H_A^\tau(h^\tau)} \mu^*(h^\tau)(v, \bar{h}_A^\tau) \pi_v^*(\widetilde{h_A^\tau}, \mathbf{M}_\tau) > 0.$

Before dealing with the zero probability events, note that the above transformation can be applied to all $(v, h_A^t)$, $h_A^t \in H_A^t(h^t)$, regardless of whether $\mu^*(h^t)(v, h_A^t) > 0$. This is because sequential rationality applies to all information sets of the agent and, thus, to all $h_A^t \in H_A^t(h^t)$. A consequence of the above transformation is that as long as the principal assigns positive probability to the event that the agent's type is $v \in V$, the agent plays the same at all of her payoff-irrelevant histories; even at those to which the principal assigns zero probability.

Now consider those types $v \in V$ such that $\sum_{h_A^t \in H_A^t(h^t)} \mu^*(h^t)(v, h_A^t) = 0$. For any

---

$h_A^t \in H_A^t(h^t)$, she participates with probability

$$\pi_v^{**}(h_A^t, \mathbf{M}_t) = \sum_{h_A^{t\,'}} \frac{\pi_v^*(h_A^{t\,'}, \mathbf{M}_t)}{|H_A^t(h_A^t)|},$$

Suppose that $(h^t, \mathbf{M}_t, 1)$ has positive probability conditional on the principal offering $\mathbf{M}_t$. Then, modify the agent's reporting strategy so that she reports $m \in M^{\mathbf{M}_t}$ with probability:

$$r_v^{**}(h_A^t, \mathbf{M}_t, 1)(m) = \sum_{h_A^{t\,'} \in H_A^t(h^t)} \frac{\pi_v^*(h_A^{t\,'}, \mathbf{M}_t)}{\sum_{\widetilde{h_A^t} \in H_A^t(h^t)} \pi_v^*(\widetilde{h_A^t}, \mathbf{M}_t)} r_v^*(h_A^{t\,'}, \mathbf{M}_t, 1)(m),$$

if $\sum_{\widetilde{h_A^t} \in H_A^t(h^t)} \pi_v^*(\widetilde{h_A^t}, \mathbf{M}_t) > 0$ and with probability

$$r_v^{**}(h_A^t, \mathbf{M}_t, 1)(m) = \sum_{h_A^{t\,'} \in H_A^t(h^t)} \frac{r_v^*(h_A^{t\,'}, \mathbf{M}_t, 1)(m)}{|H_A^t(h^t)|},$$

otherwise.

For each $h^\tau$ that has positive probability from $(h^t, \mathbf{M}_t)$ and $h_A^\tau \in H_A^\tau(h^\tau)$, for each $\mathbf{M}_\tau \in \text{supp } \Gamma^*(h^\tau)$, the agent participates with probability

$$\pi_v^{**}(h_A^\tau, \mathbf{M}_\tau) = \sum_{h_A^{\tau\,'} \in H_A^\tau(h^\tau)} \frac{\pi_v^*(h_A^{\tau\,'}, \mathbf{M}_\tau)}{|H_A^\tau(h^\tau)|}. \tag{7}$$

If the principal assigns positive probability to $(h^\tau, \mathbf{M}_\tau, 1)$ at $h^\tau$, then modify the agent's reporting strategy so that she reports $m \in M^{\mathbf{M}_\tau}$ with probability

$$r_v^{**}(h_A^\tau, \mathbf{M}_\tau, 1)(m) = \sum_{h_A^{\tau\,'} \in H_A^\tau(h^\tau)} \frac{\pi_v^*(h_A^{\tau\,'}, \mathbf{M}_\tau)}{\sum_{\bar{h}_A^\tau \in H_A^\tau(h^\tau)} \pi_v^*(\widetilde{h_A^\tau}, \mathbf{M}_\tau)} r_v^*(h_A^{\tau\,'}, \mathbf{M}_\tau, 1)(m), \tag{8}$$

if $\sum_{\bar{h}_A^\tau \in H_A^\tau(h^\tau)} \pi_v^*(\widetilde{h_A^\tau}, \mathbf{M}_\tau) > 0$ and with probability

$$r_v^{**}(h_A^\tau, \mathbf{M}_\tau, 1)(m) = \sum_{h_A^{\tau\,'} \in H_A^\tau(h^t)} \frac{r_v^*(h_A^{\tau\,'}, \mathbf{M}_\tau, 1)(m)}{|H_A^\tau(h^\tau)|},$$

otherwise. Thus, for those valuations $v \in V$ to which the principal assigns 0 probability–either at $h^t$ or at a continuation history $h^\tau$ on the equilibrium path of the strategy given $h^{t-}$ their strategies also do not depend on the payoff-irrelevant part of the private history.

Fix $\tau \geqslant t$. Under the new strategy, the principal's beliefs that the agent is of type $v$ and her private history is $h_A^{\tau+1}$ at history $h^{\tau+1} = (h^\tau, \mathbf{M}_\tau, 1, s_\tau, a_\tau)$, $\mathbf{M}_\tau \in \text{supp } \Gamma^*(h^\tau)$ are given by:

$$\mu^{**}(h^{\tau+1})(v, h_A^{\tau+1}) = \frac{\mu^{**}(h^\tau)(v, h_A^\tau)\pi_v^*(h_A^\tau, \mathbf{M}_\tau)r_v^{**}(h_A^\tau, \mathbf{M}_\tau, 1)(m_\tau)\beta^{\mathbf{M}_\tau}(s_\tau|m_\tau)}{\sum_{\tilde{v},\widetilde{h_A^\tau}} \sum_{\tilde{m}_\tau} \mu^{**}(\tilde{v}, \widetilde{h_A^\tau})\pi_{\tilde{v}}^{**}(\widetilde{h_A^\tau}, \mathbf{M}_\tau)r_{\tilde{v}}^{**}(\widetilde{h_A^\tau}, \mathbf{M}_\tau, 1)(\tilde{m}_\tau)\beta^{\mathbf{M}_\tau}(s_\tau|\tilde{m}_\tau)}, \tag{9}$$

where $\mu^{**}(h^t)(v, h_A^t) = \mu^*(h^t)(v, h_A^t)$ and at history $h^{\tau+1} = (h^\tau, \mathbf{M}_\tau, 0, s_\tau, a_\tau)$, $\mathbf{M}_\tau \in$ supp $\Gamma^*(h^\tau)$ are given by:

$$\mu^{**}(h^{\tau+1})(v, h_A^{\tau+1}) = \frac{\mu^{**}(h^\tau)(v, h_A^\tau)(1 - \pi_v^{**}(h_A^\tau, \mathbf{M}_\tau))}{\sum_{\tilde{v},\widetilde{h_A^\tau}} \mu^{**}(h^\tau)(\tilde{v}, \widetilde{h_A^\tau})(1 - \pi_{\tilde{v}}^{**}(\widetilde{h_A^\tau}, \mathbf{M}_\tau))}, \tag{10}$$

We now show by induction that for any $\tau \geqslant t$,

$$\sum_{h_A^{\tau+1} \in H_A^{\tau+1}(h^{\tau+1})} \mu^{**}(h^{\tau+1})(v, h_A^{\tau+1}) = \sum_{h_A^{\tau+1} \in H_A^{\tau+1}(h^{\tau+1})} \mu^*(h^{\tau+1})(v, h_A^{\tau+1}).$$

We do so for histories $h^\tau$ that are consistent with the equilibrium strategy for which the agent participates; it is immediate that the same holds for those histories in which she does not.

For $\tau = t$ and $h^{t+1} = (h^t, \mathbf{M}_t, 1, s_t, a_t)$, note the denominator on the right-hand side of equation (9) can be written as:

$$\sum_{\tilde{v},h_A^t} \sum_{m \in M^{\mathbf{M}_t}} \mu^*(h^t)(\tilde{v}, h_A^t)\pi_{\tilde{v}}^{**}(h_A^t, \mathbf{M}_t)r_{\tilde{v}}^{**}(h_A^t, \mathbf{M}_t, 1)(m)\beta^{\mathbf{M}_t}(s|m)$$

$$= \sum_{\tilde{v}} \sum_{m \in M^{\mathbf{M}_t}} \pi_{\tilde{v}}^{**}(h_A^t, \mathbf{M}_t)r_{\tilde{v}}^{**}(h_A^t, \mathbf{M}_t, 1)(m)\beta^{\mathbf{M}_t}(s|m) \sum_{h_A^t} \mu^*(h^t)(\tilde{v}, h_A^t)$$

$$= \sum_{\tilde{v},h_A^t} \sum_{m \in M^{\mathbf{M}_t}} \mu^*(h^t)(\tilde{v}, h_A^t)\pi_{\tilde{v}}^*(h_A^t, \mathbf{M}_t)r_{\tilde{v}}^*(h_A^t, \mathbf{M}_t, 1)(m)\beta^{\mathbf{M}_t}(s|m),$$

where the first equality uses that $\pi_v^{**}, r_v^{**}$ does not depend on $h_A^t$ and the second equality uses the definition of $\pi_v^{**}, r_v^{**}$; see equations (7) and (8). Note the last expression is the denominator in $\mu^*(h^{t+1})(v, h_A^{t+1})$. Therefore, for $h^{t+1} = (h^t, \mathbf{M}_t, 1, s_t, a_t), h_A^{t+1} = (h_A^t, \mathbf{M}_t, 1, m_t, s_t, a_t), h_A^t \in H_A^t(h^t)$

$$
\begin{aligned}
\sum_{h_A^t} \mu^{**}(h^{t+1})(v, h_A^{t+1}) &= \frac{\sum_{h_A^t} \mu^*(h^t)(v, h_A^t) \pi_v^{**}(h_A^t, \mathbf{M}_t) r_v^{**}(h_A^t, \mathbf{M}_t, 1)(m_t) \beta^{\mathbf{M}_t}(s_t|m_t)}{\sum_{\tilde{v}, \widetilde{h_A^t}} \sum_{m \in M^{\mathbf{M}_t}} \mu^*(h^t)(\tilde{v}, \widetilde{h_A^t}) \pi_{\tilde{v}}^{**}(\widetilde{h_A^t}, \mathbf{M}_t) r_{\tilde{v}}^{**}(\widetilde{h_A^t}, \mathbf{M}_t, 1)(m) \beta^{\mathbf{M}_t}(s_t|m)} \\
&= \frac{\pi_v^{**}(h_A^t, \mathbf{M}_t) r_v^{**}(h_A^t, \mathbf{M}_t, 1)(m_t) \beta^{\mathbf{M}_t}(s_t|m_t) \sum_{h_A^t} \mu^*(h^t)(v, h_A^t)}{\sum_{\tilde{v}, \widetilde{h_A^t}} \sum_{m \in M^{\mathbf{M}_t}} \mu^*(h^t)(\tilde{v}, \widetilde{h_A^t}) \pi_{\tilde{v}}^{**}(\widetilde{h_A^t}, \mathbf{M}_t, 1) r_{\tilde{v}}^{**}(\widetilde{h_A^t}, \mathbf{M}_t, 1)(m) \beta^{\mathbf{M}_t}(s_t|m)} \\
&= \frac{\sum_{h_A^t} \mu^*(h^t)(v, h_A^t) \pi_v^*(h_A^t, \mathbf{M}_t) r_v^*(h_A^t, \mathbf{M}_t, 1)(m_t) \beta^{\mathbf{M}_t}(s_t|m_t)}{\sum_{\tilde{v}, \widetilde{h_A^t}} \sum_{m \in M^{\mathbf{M}_t}} \mu^*(h^t)(\tilde{v}, \widetilde{h_A^t}) \pi_{\tilde{v}}^{**}(\widetilde{h_A^t}, \mathbf{M}_t) r_{\tilde{v}}^{**}(\widetilde{h_A^t}, \mathbf{M}_t, 1)(m) \beta^{\mathbf{M}_t}(s_t|m)} \\
&= \frac{\sum_{h_A^t} \mu^*(h^t)(v, h_A^t) \pi_v^*(h_A^t, \mathbf{M}_t) r_v^*(h_A^t, \mathbf{M}_t, 1)(m_t) \beta^{\mathbf{M}_t}(s_t|m_t)}{\sum_{\tilde{v}, \widetilde{h_A^t}} \sum_{m \in M^{\mathbf{M}_t}} \mu^*(h^t)(\tilde{v}, \widetilde{h_A^t}) \pi_{\tilde{v}}^*(\widetilde{h_A^t}, \mathbf{M}_t) r_{\tilde{v}}^*(\widetilde{h_A^t}, \mathbf{M}_t, 1)(m) \beta^{\mathbf{M}_t}(s_t|m)} \\
&= \sum_{h_A^t} \mu^*(h^{t+1})(v, h_A^{t+1}),
\end{aligned}
$$

where the second equality uses that $\pi_v^{**}, r_v^{**}$ do not depend on $h_A^t$, the third equality uses the definition of $\pi_v^{**}, r_v^{**}$, the fourth equality uses the observation about the denominator, and the last equality follows by definition of $\mu^*(h^{t+1})$. Adding up both sides of the expression over $m_t$ delivers the desired expression.

Now suppose we have established the above claim for each $\tau' < \tau$. We now show that it holds to $\tau' = \tau$. To see that it holds for $\tau' = \tau$, note the probability of $h^{\tau+1} = (h^\tau, \mathbf{M}_\tau, 1, s_\tau)$ conditional on $h^\tau$, which is given by the denominator on the right hand

side of equation (9), can be written as:

$$\sum_{\tilde{v},\widetilde{h_A^\tau}} \sum_{\tilde{m}_\tau} \mu^{**}(\tilde{v}, \widetilde{h_A^\tau}) \pi_{\tilde{v}}^{**}(\widetilde{h_A^\tau}, \mathbf{M}_\tau) r_{\tilde{v}}^{**}(\widetilde{h_A^\tau}, \mathbf{M}_\tau, 1)(\tilde{m}_\tau) \beta^{\mathbf{M}_\tau}(s_\tau | \tilde{m}_\tau)$$

$$= \sum_{\tilde{v}} \sum_{\tilde{m}_\tau} \pi_{\tilde{v}}^{**}(\widetilde{h_A^\tau}, \mathbf{M}_\tau) r_v^{**}(\widetilde{h_A^\tau}, \mathbf{M}_\tau, 1)(\tilde{m}_\tau) \sum_{\widetilde{h_A^\tau}} \mu^{**}(\tilde{v}, \widetilde{h_A^\tau}) \beta^{\mathbf{M}_\tau}(s_\tau | \tilde{m}_\tau)$$

$$= \sum_{\tilde{v}} \sum_{\tilde{m}_\tau} \pi_v^{**}(\widetilde{h_A^\tau}, \mathbf{M}_\tau) r_v^{**}(\widetilde{h_A^\tau}, \mathbf{M}_\tau, 1)(\tilde{m}_\tau) \sum_{\widetilde{h_A^\tau}} \mu^{*}(\tilde{v}, \widetilde{h_A^\tau}) \beta^{\mathbf{M}_\tau}(s_\tau | \tilde{m}_\tau)$$

$$= \sum_{\tilde{v},\widetilde{h_A^\tau}} \sum_{\tilde{m}_\tau} \mu^{*}(\tilde{v}, \widetilde{h_A^\tau}) \pi_v^{*}(\widetilde{h_A^\tau}, \mathbf{M}_\tau) r_v^{*}(\widetilde{h_A^\tau}, \mathbf{M}_\tau, 1)(\tilde{m}_\tau) \beta^{\mathbf{M}_\tau}(s_\tau | \tilde{m}_\tau),$$

where the second equality makes use of the inductive hypothesis for $\tau' = \tau-1$, $\sum_{\widetilde{h_A^\tau}} \mu^{**}(h^\tau)(\tilde{v}, \widetilde{h_A^\tau}) = \sum_{\widetilde{h_A^\tau}} \mu^{*}(h^\tau)(\tilde{v}, \widetilde{h_A^\tau})$, and the third equality uses the definition of the participation and reporting strategies defined in equations (7) and (8). Note the last line corresponds to the expression of the denominator of $\mu^{*}(h^{\tau+1})(v, h_A^{\tau+1})$ in the original PBE.

Therefore,

$$\sum_{h_A^\tau} \mu^{**}(h^{\tau+1})(v, h_A^{\tau+1}) = \frac{\sum_{h_A^\tau} \mu^{**}(h^\tau)(v, h_A^\tau) \pi_v^{**}(h_A^\tau, \mathbf{M}_\tau) r_v^{**}(h_A^\tau, \mathbf{M}_\tau, 1)(m_\tau) \beta^{\mathbf{M}_\tau}(s_\tau | m_\tau)}{\sum_{\tilde{v},\widetilde{h_A^\tau}} \sum_{\tilde{m}_\tau} \mu^{*}(h^\tau)(\tilde{v}, \widetilde{h_A^\tau}) \pi_{\tilde{v}}^{*}(\widetilde{h_A^\tau}, \mathbf{M}_\tau) r_{\tilde{v}}^{*}(\widetilde{h_A^\tau}, \mathbf{M}_\tau, 1)(\tilde{m}_\tau) \beta^{\mathbf{M}_\tau}(s_\tau | \tilde{m}_\tau)}$$

$$= \frac{\pi_v^{**}(h_A^\tau, \mathbf{M}_\tau) r_v^{**}(h_A^\tau, \mathbf{M}_\tau, 1)(m_\tau) \beta^{\mathbf{M}_\tau}(s_\tau | m_\tau) \sum_{h_A^\tau} \mu^{**}(h^\tau)(v, h_A^\tau)}{\sum_{\tilde{v},\widetilde{h_A^\tau}} \sum_{\tilde{m}_\tau} \mu^{*}(h^\tau)(\tilde{v}, \widetilde{h_A^\tau}) \pi_{\tilde{v}}^{*}(\widetilde{h_A^\tau}, \mathbf{M}_\tau) r_{\tilde{v}}^{*}(\widetilde{h_A^\tau}, \mathbf{M}_\tau, 1)(\tilde{m}_\tau) \beta^{\mathbf{M}_\tau}(s_\tau | \tilde{m}_\tau)}$$

$$= \frac{\pi_v^{**}(h_A^\tau, \mathbf{M}_\tau) r_v^{**}(h_A^\tau, \mathbf{M}_\tau, 1)(m_\tau) \beta^{\mathbf{M}_\tau}(s_\tau | m_\tau) \sum_{h_A^\tau} \mu^{*}(h^\tau)(v, h_A^\tau)}{\sum_{\tilde{v},\widetilde{h_A^\tau}} \sum_{\tilde{m}_\tau} \mu^{*}(h^\tau)(\tilde{v}, \widetilde{h_A^\tau}) \pi_{\tilde{v}}^{*}(\widetilde{h_A^\tau}, \mathbf{M}_\tau) r_{\tilde{v}}^{*}(\widetilde{h_A^\tau}, \mathbf{M}_\tau, 1)(\tilde{m}_\tau) \beta^{\mathbf{M}_\tau}(s_\tau | \tilde{m}_\tau)}$$

$$= \frac{\sum_{h_A^\tau} \mu^{*}(h^\tau)(v, h_A^\tau) \pi_v^{*}(h_A^\tau, \mathbf{M}_\tau) r_v^{*}(h_A^\tau, \mathbf{M}_\tau, 1)(m_\tau) \beta^{\mathbf{M}_\tau}(s_\tau | m_\tau)}{\sum_{\tilde{v},\widetilde{h_A^\tau}} \sum_{\tilde{m}_\tau} \mu^{*}(h^\tau)(\tilde{v}, \widetilde{h_A^\tau}) \pi_{\tilde{v}}^{*}(\widetilde{h_A^\tau}, \mathbf{M}_\tau) r_{\tilde{v}}^{*}(\widetilde{h_A^\tau}, \mathbf{M}_\tau, 1)(\tilde{m}_\tau) \beta^{\mathbf{M}_\tau}(s_\tau | \tilde{m}_\tau)}$$

$$= \sum_{h_A^\tau} \mu^{*}(h^{\tau+1})(v, h_A^{\tau+1}),$$

where the first equality makes use of our conclusion for the denominator, the second equality uses that $\pi_v^{**}, r_v^{**}$ do not depend on $h_A^\tau$, the third equality makes use of the inductive hypothesis, the fourth equality makes use of equations (7) and (8), and the fifth equality follows from the definition of the beliefs via Bayes' rule. Adding up the

above expression over $m_\tau$ delivers the desired conclusion.

We now use the above properties to show the payoffs of the principal do not change along the path that starts from $h^t$. Fix any history $h^\tau$ that is on the path of the equilibrium strategy starting from $h^t$. The principal's payoffs are given by:

$$\sum_{(v,h_A^\tau)} \mu^{**}(h^\tau)(v,h_A^\tau) \left\{ (1 - \pi_v^{**}(v,h_A^\tau)) \mathbb{E}^{\Gamma^*,\pi_v^{**},r_v^{**}} \left[ W(a(h^\tau),a^*,a^{\geq\tau+1})|v,h_A^\tau,\mathbf{M}_\tau,0 \right] + \pi_v^{**}(h_A^\tau,\mathbf{M}_\tau) \times \right.$$

$$\left. \sum_{m\in M^{\mathbf{M}_\tau}} r_v^{**}(h_A^\tau,\mathbf{M}_\tau,1)(m) \sum_{s\in S^{\mathbf{M}_\tau}} \beta^{\mathbf{M}_\tau}(s|m) \mathbb{E}^{\Gamma^*,\pi_v^{**},r_v^{**}} \left[ W(a(h^\tau),a_s,a^{\geq\tau+1},v)|v,h_A^\tau,1,m,s \right] \right\}$$

$$= \sum_v \left\{ (1 - \pi_v^{**}(h_A^\tau,\mathbf{M}_\tau)) \sum_{h_A^\tau} \mu^{**}(h^\tau)(v,h_A^\tau) \mathbb{E}^{\Gamma^*,\pi_v^{**},r_v^{**}} \left[ W(a(h^\tau),a^*,a^{\geq\tau+1})|v,h_A^\tau,\mathbf{M}_\tau,0 \right] + \pi_v^{**}(h_A^\tau,\mathbf{M}_\tau) \times \right.$$

$$\left. \sum_{m\in M^{\mathbf{M}_\tau}} r_v^{**}(h_A^\tau,\mathbf{M}_\tau,1)(m) \sum_{h_A^\tau} \mu^{**}(h^\tau)(v,h_A^\tau) \sum_{s\in S^{\mathbf{M}_\tau}} \beta^{\mathbf{M}_\tau}(s|m) \mathbb{E}^{\Gamma^*,\pi_v^{**},r_v^{**}} \left[ W(a(h^\tau),a_s,a^{\geq\tau+1},v)|v,h_A^\tau,1,m,s \right] \right\}$$

$$= \sum_v \left\{ (1 - \pi_v^{*}(h_A^\tau,\mathbf{M}_\tau)) \sum_{h_A^\tau} \mu^{*}(h^\tau)(v,h_A^\tau) \mathbb{E}^{\Gamma^*,\pi_v^{**},r_v^{**}} \left[ W(a(h^\tau),a^*,a^{\geq\tau+1})|v,h_A^\tau,\mathbf{M}_\tau,0 \right] + \pi_v^{**}(h_A^\tau,\mathbf{M}_\tau) \times \right.$$

$$\left. \sum_{m\in M^{\mathbf{M}_\tau}} r_v^{**}(h_A^\tau,\mathbf{M}_\tau,1)(m) \sum_{h_A^\tau} \mu^{*}(h^\tau)(v,h_A^\tau) \sum_{s\in S^{\mathbf{M}_\tau}} \beta^{\mathbf{M}_\tau}(s|m) \mathbb{E}^{\Gamma^*,\pi_v^{**},r_v^{**}} \left[ W(a(h^\tau),a_s,a^{\geq\tau+1},v)|v,h_A^\tau,1,m,s \right] \right\}$$

$$= \sum_{(v,h_A^\tau)} \mu^{*}(h^\tau)(v,h_A^\tau) \left\{ (1 - \pi_v^{*}(v,h_A^\tau)) \mathbb{E}^{\Gamma^*,\pi_v^{**},r_v^{**}} \left[ W(a(h^\tau),a^*,a^{\geq\tau+1})|v,h_A^\tau,\mathbf{M}_\tau,0 \right] \right.$$

$$+ \pi_v^{*}(h_A^\tau,\mathbf{M}_\tau) \times \sum_{m\in M^{\mathbf{M}_\tau}} r_v^{*}(h_A^\tau,\mathbf{M}_\tau,1)(m) \sum_{s\in S^{\mathbf{M}_\tau}} \beta^{\mathbf{M}_\tau}(s|m) \mathbb{E}^{\Gamma^*,\pi_v^{**},r_v^{**}} \left[ W(a(h^\tau),a_s,a^{\geq\tau+1},v)|v,h_A^\tau,1,m,s \right] \right\}$$

$$= \sum_{(v,h_A^\tau)} \mu^{*}(h^\tau)(v,h_A^\tau) \left\{ (1 - \pi_v^{*}(v,h_A^\tau)) \mathbb{E}^{\Gamma^*,\pi_v^{*},r_v^{*}} \left[ W(a(h^\tau),a^*,a^{\geq\tau+1})|v,h_A^\tau,\mathbf{M}_\tau,0 \right] \right.$$

$$+ \pi_v^{*}(h_A^\tau,\mathbf{M}_\tau) \times \sum_{m\in M^{\mathbf{M}_\tau}} r_v^{*}(h_A^\tau,\mathbf{M}_\tau,1)(m) \sum_{s\in S^{\mathbf{M}_\tau}} \beta^{\mathbf{M}_\tau}(s|m) \mathbb{E}^{\Gamma^*,\pi_v^{*},r_v^{*}} \left[ W(a(h^\tau),a_s,a^{\geq\tau+1},v)|v,h_A^\tau,1,m,s \right] \right\},$$

where the first equality follows from noting that the agent's strategy does not depend on $h_A^\tau$, the second equality follows from the previous result and noting that under $\pi_v^{**}, r_v^{**}$, the continuation strategy does not depend on $h_A^\tau$, the third equality follows from the definition of the strategy, and the last equality follows from noting this equality holds after every history on the path of the equilibrium strategy starting from $h^t$. $\square$

49

**Remark A.1.** Suppose that at history $h^t$ and after offering $\mathbf{M}_t$, the principal assigns probability 0 to the agent participating. In that case, his beliefs $\mu^*(h^t, \mathbf{M}_t, 1) \in \Delta(V \times H_A^t(h^t, \mathbf{M}_t, 1))$ are not determined by Bayes' rule. It is immediate to extend the proof of Proposition A.1 to show that starting from $(h^t, \mathbf{M}_t, 1)$, we can modify the agent's strategy along the path of the equilibrium strategy so that she does not condition on her payoff-irrelevant private history.

**Remark A.2.** The payoff-equivalent PBE assessment one obtains from Proposition A.1 satisfies the following property. On the equilibrium path, the principal's beliefs over the agent's payoff-relevant type, $v \in V$, do not depend on her payoff-irrelevant history, $h_A^t$. That is, for any public history on the equilibrium path of the strategy profile given the initial history, for any $v \in V$, $h_A^t, h_A^{t\,\prime} \in H_A^t(h^t)$ such that $\mu^*(h^t)(v, h_A^t), \mu^*(h^t)(v, h_A^{t\,\prime}) > 0$, we have $\mu^*(h^t)(v, h_A^t) = \mu^*(h^t)(v, h_A^{t\,\prime})$.

However, at a public history $h^t$ reached after a deviation by the agent, either because she changes her participation strategy in a detectable way or because she triggers an output message that was not supposed to be triggered according to the equilibrium strategy, the requirements of PBE do not rule out that the principal's *updated* beliefs depend non-trivially on both $v$ and $h_A^t$.

However, it follows from Proposition A.1 that without loss of generality, we can assume that when the principal observes a deviation by the agent, his updated beliefs do not depend on $h_A^t$. After all, the agent's behavior after the deviation does not depend on his payoff irrelevant private history and the principal cannot offer mechanisms as a function of $h_A^t$. We record this in Corollary A.1 below and prove it in Section III of the supplementary material.

**Corollary A.1.** *Fix a PBE, $\langle \Gamma^*, (\pi_v^*, r_v^*)_{v \in V}, \mu^* \rangle$. Then, without loss of generality, for any $t$, public history $h^t$, $v \in V$ and private histories $h_A^t, h_A^{t\,\prime} \in H_A^t(h^t)$ such that $\mu^*(h^t)(v, h_A^t), \mu^*(h^t)(v, h_A^{t\,\prime}) > 0$,*

$$\mu^*(h^t)(v, h_A^t) = \mu^*(h^t)(v, h_A^{t\,\prime}).$$

**Remark A.3.** Note that the corollary states that equality holds only for $(v, h_A^t), (v, h_A^{t\,\prime})$ that have positive probability given the equilibrium beliefs, because the agent's strategy

may assign probability 0 to some input messages and hence some $h_A^t$'s cannot be given positive probability.

Hereafter, we focus on equilibria in which the agent's strategy does not depend on the payoff-irrelevant part of her private history.

We introduce one final piece of notation. Given a strategy profile $(\Gamma^*, (\pi_v^*, r_v^*)_{v \in V})$ and a history $h^t$, denote the continuation strategy starting at $h^t$ implied by this profile as $(\Gamma^*, (\pi_v^*, r_v^*)_{v \in V})\big|_{h^t}$.

The next result shows that without loss of generality, we can focus on equilibria in which the agent participates with probability 1:

**Proposition A.2.** *Fix a PBE,* $\langle \Gamma^*, (\pi_v^*, r_v^*)_{v \in V}, \mu^* \rangle$. *Then, there is a PBE,* $\langle \Gamma^{**}, (\pi_v^{**}, r_v^{**})_{v \in V}, \mu^{**} \rangle$ *where*

1. *For every* $t \geqslant 0$, *for every* $v \in V$, *for every* $h_A^t$, $\pi_v^{**}(h_A^t, \mathbf{M}_t) = 1$ *for all* $\mathbf{M}_t \in$ *supp* $\Gamma^{**}(h^t)$.

2. *The principal and the agent's payoffs are the same after every history as in* $\langle \Gamma^*, (\pi_v^*, r_v^*)_{v \in V}, \mu^* \rangle$.

3. *For every* $t$ *and public history* $h^t$, *the distribution over allocations is the same as in* $\langle \Gamma^*, (\pi_v^*, r_v^*)_{v \in V}, \mu^* \rangle$.

*Proof.* Fix $t$ and $h^t$ such that there exists $\mathbf{M}_t \in$ supp $\Gamma^*(h^t)$ such that $\pi_v^*(h_A^t, \mathbf{M}_t) < 1$. Recall $M^{\mathbf{M}_t}$ is a finite set and $S^{\mathbf{M}_t}$ contains $\Delta(M^{\mathbf{M}_t})$. Recall that for all $m \in M^{\mathbf{M}_t}$, $\beta^{\mathbf{M}_t}(\cdot|m) \in \Delta^*(S^{\mathbf{M}_t})$ has finite support. Then there exists $s^* \in S^{\mathbf{M}_t}$ such that $\beta^{\mathbf{M}_t}(s^*|m) = 0$ for all $m \in M^{\mathbf{M}_t}$.

Let $V_1 = \{v \in V : \pi_v^*(h_A^t, \mathbf{M}_t) > 0\} = \{v_1, ..., v_{|V_1|}\}$. Since $|V| \leqslant |M^{\mathbf{M}_t}|$, we can label the latter set $M^{\mathbf{M}_t} = \{m_1, ...., m_{|V_1|}, ..., m_{|M^{\mathbf{M}_t}|}\}$. Modify $\beta^{\mathbf{M}_t}$ as follows. For $i = 1, \ldots, |V_1|$, let

$$\widetilde{\beta}(s|m_i) = \sum_{m \in M^{\mathbf{M}_t}} \beta^{\mathbf{M}_t}(s|m) r_{v_i}^*(h_A^t, \mathbf{M}_t, 1)(m).$$

Note that $\widetilde{\beta}$ does not depend on $h_A^t$ since $r_{v_i}^*$ does not depend on $h_A^t$.

If $|V_1| < |M^{\mathbf{M}_t}|$, let $\widetilde{\beta}(s^*|m_i) = 1$ for all $i > |V_1|$ and let $\tilde{\alpha}(s^*) = \delta_{a^*}$. Modify the

strategies so that the principal, instead of offering $\mathbf{M}_t$, offers $\widetilde{\mathbf{M}}_t = \{\langle M^{\mathbf{M}_t}, \tilde{\beta}, S^{\mathbf{M}_t}\rangle, \tilde{\alpha}\}$

$$r^{**}_{v_i}(h^t_A, \mathbf{M}_t, 1)(m) = \begin{cases} \pi^*_{v_i}(h^t_A, \mathbf{M}_t) & \text{if } m = m_i \\ (1 - \pi^*_{v_i}(h^t_A, \mathbf{M}_t)) & \text{if } m = m_{|V_1|+1} \\ 0 & \text{otherwise} \end{cases}$$

and let $(\Gamma^{**}, (\pi^{**}_v, r^{**}_v)_{v\in V})\big|_{(h^t, \tilde{\mathbf{M}}_t, 1, s^*)} = (\Gamma^*, (\pi^*_v, r^*_v)_{v\in V})\big|_{(h^t, \mathbf{M}_t, 0)}$, for all other $s \in S^{\mathbf{M}_t}$, let $(\Gamma^{**}, (\pi^{**}_v, r^{**}_v)_{v\in V})\big|_{(h^t, \tilde{\mathbf{M}}_t, 1, s)} = (\Gamma^*, (\pi^*_v, r^*_v)_{v\in V})\big|_{(h^t, \mathbf{M}_t, 1, s)}$.

If $|V_1| = |M^{\mathbf{M}_t}|$ (which implies that $V_1 = V$), modify $\tilde{\beta}$ once more so that:

$$\tilde{\tilde{\beta}}(s|m_i) = \begin{cases} \pi^*_{v_i}(h^t_A, \mathbf{M}_t)\tilde{\beta}(s|m_i) & \text{if } s \neq s^* \\ (1 - \pi^*_{v_i}(h^t_A, \mathbf{M}_t)) & \text{otherwise} \end{cases},$$

and let $\tilde{\alpha}(s^*) = \delta_{a^*}$ as before. Modify the strategies so that the principal, instead of offering $\mathbf{M}_t$, offers $\widetilde{\mathbf{M}}_t = \{\langle M^{\mathbf{M}_t}, \tilde{\tilde{\beta}}, S^{\mathbf{M}_t}\rangle, \tilde{\alpha}\}$

$$r^{**}_{v_i}(h^t_A, \mathbf{M}_t, 1)(m) = \mathbb{1}[m = m_i],$$

and let $(\Gamma^{**}, (\pi^{**}_v, r^{**}_v)_{v\in V})\big|_{(h^t, \tilde{\mathbf{M}}_t, 1, s^*)} = (\Gamma^*, (\pi^*_v, r^*_v)_{v\in V})\big|_{(h^t, \mathbf{M}_t, 0)}$, for all other $s \in S^{\mathbf{M}_t}$, let $(\Gamma^{**}, (\pi^{**}_v, r^{**}_v)_{v\in V})\big|_{(h^t, \tilde{\mathbf{M}}_t, 1, s)} = (\Gamma^*, (\pi^*_v, r^*_v)_{v\in V})\big|_{(h^t, \mathbf{M}_t, 1, s)}$.

It follows immediately that the principal's and the agent's payoffs remain the same and we have not changed the distribution over allocations at any history starting from $h^t$. $\square$

**Remark A.4.** To keep things simple, the proof of Proposition A.2 uses the restriction that $\beta^{\mathbf{M}_t}$ has finite support and $S^{\mathbf{M}_t}$ is a *large* set to add an output message that allows the principal to

1. replicate the agent's non-participation decision and,

2. make it incentive compatible for him to offer the same continuation upon observing $a^*$ as he was offering before.

One can write an albeit more notationally involved proof that (i) does not rely on the existence of an output message that is never sent and (ii) respects the one-to-one mapping between posteriors and output messages. This alternative proof is available from the authors upon request.

**Proposition A.3.** *Fix a PBE, $\langle \Gamma^*, (\pi_v^*, r_v^*)_{v \in V}, \mu^* \rangle$ that satisfies the properties of Propositions A.1 and A.2. Then, without loss of generality, there is a one-to-one map between output messages and continuation beliefs. That is, for every $t$, public history $h^t$, $\mathbf{M}_t \in$ supp $\Gamma^*(h^t)$, if $s_t, s_t' \in S^{*\mathbf{M}_t}$ is such that $s_t \neq s_t'$, then $\mu^*(h^t, \mathbf{M}_t, 1, s_t) \neq \mu^*(h^t, \mathbf{M}_t, 1, s_t')$.*

Lemma A.1 is used to prove Proposition A.3:

**Lemma A.1.** Fix a PBE assessment, $\langle \Gamma^*, (\pi_v^*, r_v^*)_{v \in V}, \mu^* \rangle$, that satisfies the properties of Proposition A.1 and A.2. Then, there is another assessment $\langle \Gamma^{**}, (\pi_v^{**}, r_v^{**})_{v \in V}, \mu^{**} \rangle$ that satisfies the properties of Proposition A.1 and the following holds:

1. For all $h^t$, for all $\mathbf{M}_t \in$ supp $\Gamma^{**}(h^t)$, $\pi_v^{**}(h_A^t, \mathbf{M}_t) = 1$ for all $v \in V$ such that $\sum_{h_A^t} \mu^*(h^t)(v, h_A^t) > 0$.

2. For all $h^t$, for all $\mathbf{M}_t \in$ supp $\Gamma^{**}(h^t)$, if $s \in S^{*\mathbf{M}_t}$, then

$$\sum_{(v, h_A^t), m \in M^{\mathbf{M}_t}} \mu^{**}(h^t)(v, h_A^t) r_v^{**}(h_A^t, \mathbf{M}_t)(m) \beta^{\mathbf{M}_t}(s|m) > 0.$$

3. For all $h^t$, the principal's continuation payoff remains the same and he faces the same distribution over allocations at each continuation history on the path of the equilibrium strategy given $h^t$. The same holds for the agent for each of her types $v \in V$ which have positive probability at $h^t$.

Among other things, Lemma A.1 guarantees that if the principal's strategy specifies that mechanism $\mathbf{M}_t$ is played at history $h^t$, then any output message $s \in S^{*\mathbf{M}_t} (\equiv \{s \in S^{\mathbf{M}_t} : (\exists m \in M^{\mathbf{M}_t}) \beta^{\mathbf{M}_t}(s|m) > 0\})$ has positive probability under the equilibrium strategy profile. Thus, the principal is never *surprised* by the output messages he observes.

*Proof of Lemma A.1.* Consider a PBE assessment, $\langle \Gamma^*, (\pi_v^*, r_v^*)_{v \in V}, \mu^* \rangle$, that satisfies the properties of Propositions A.1 and A.2. Suppose there exists a history $h^t$ and a type $v \in V$ to which the principal assigns probability 0. That is, $\sum_{h_A^t \in H_A^t(h_A^t)} \mu^*(h^t)(v, h_A^t) = 0$. Let $\mathbf{M}_t \in$ supp $\Gamma^*(h^t)$. Let $M^{+\mathbf{M}_t} = \{m \in M^{\mathbf{M}_t} : \sum_{(\tilde{v}, \widetilde{h_A^t}) \in V \times H_A^t(h^t)} \mu^*(h^t)(\tilde{v}, \widetilde{h_A^t}) r_{\tilde{v}}(\widetilde{h_A^t}, \mathbf{M}_t, 1)(m) > 0\}$. If $M^{\mathbf{M}_t} \backslash M^{+\mathbf{M}_t} \neq \varnothing$, note that we can do the following transformation without upsetting the equilibrium:

First, replace $\mathbf{M}_t$ by $\mathbf{M}_t' = (\langle \beta^{\mathbf{M}_t'}, M^{\mathbf{M}_t'}, S^{\mathbf{M}_t'} \rangle, \alpha^{\mathbf{M}_t'})$ where $M^{\mathbf{M}_t'} = M^{\mathbf{M}_t}$, $S^{\mathbf{M}_t'} = $

$S^{\mathbf{M}_t}, \alpha^{\mathbf{M}'_t} = \alpha^{\mathbf{M}_t}$ and $\beta^{\mathbf{M}'_t}(\cdot|m) = \beta^{\mathbf{M}_t}(\cdot|m)$ for $m \in M^{+\mathbf{M}_t}$ and otherwise, let $\beta^{\mathbf{M}'_t}(\cdot|m) = \beta^{\mathbf{M}_t}(\cdot|m^+)$ for some $m^+ \in M^{+\mathbf{M}_t}$. Modify the principal's strategy at $h^t$ so that instead of offering $\mathbf{M}_t$, he offers $\mathbf{M}'_t$; that is, let $\Gamma^{**}(h^t)(\mathbf{M}'_t) = \Gamma^*(h^t)(\mathbf{M}_t)$

Second, modify continuation strategies so that

$$(\Gamma^{**}, (\pi_v^{**}, r_v^{**})_{v \in V})|_{(h^t, \mathbf{M}'_t, 1, s_t, a_t)} = (\Gamma^*, (\pi_v^*, r_v^*)_{v \in V})|_{(h^t, \mathbf{M}_t, 1, s_t, a_t)}$$

for those output messages and allocations consistent with $\mathbf{M}'_t$.

Third, modify the agent's strategy as follows. For $v' \in V$ such that $\sum_{h_A^t \in H_A^t(h^t)} \mu^*(h^t)(v', h_A^t) > 0$, let $\pi_{v'}^{**}(h_A^t, \mathbf{M}') = \pi_{v'}^*(h_A^t, \mathbf{M}_t) = 1$ and $r_{v'}^{**}(h_A^t, \mathbf{M}'_t, 1) = r_{v'}^*(h_A^t, \mathbf{M}_t, 1)$. Note that for these types we have not really modified the mechanism–in effect, we have removed the choices they were not making and, hence, removed possible deviations for them.

Consider now $v \in V$ such that $\sum_{h_A^t \in H_A^t(h^t)} \mu^*(h^t)(v, h_A^t) = 0$. Set $r_v^{**}(h_A^t, \mathbf{M}'_t, 1)$ so that $r_v^{**}(h_A^t, \mathbf{M}'_t, 1)(m) > 0$ if and only if $m$ solves

$$\max_{m \in M^{+\mathbf{M}_t}} \sum_{s \in S^{\mathbf{M}_t}} \beta^{\mathbf{M}_t}(s|m) \sum_{a \in A} \alpha^{\mathbf{M}_t}(m|s) \mathbb{E}^{\Gamma^*, \pi^*, r^*}[U(a(h^t), a, a^{\geq t+1}, v)|h_A^t, \mathbf{M}_t, 1, m, s, a],$$

where we are using that $\mathbf{M}'_t$ is only relabeling the input messages and the continuation histories remain the same as before so that for all $m \in M_t^{+\mathbf{M}_t}$:

$$\sum_{s \in S^{\mathbf{M}'_t}} \beta^{\mathbf{M}'_t}(s|m) \sum_{a \in A} \alpha^{\mathbf{M}'_t}(m|s) \mathbb{E}^{\Gamma^{**}, \pi^{**}, r^{**}}[U(a(h^t), a, a^{\geq t+1}, v)|h_A^t, \mathbf{M}'_t, 1, m, s, a]$$
$$= \sum_{s \in S^{\mathbf{M}_t}} \beta^{\mathbf{M}_t}(s|m) \sum_{a \in A} \alpha^{\mathbf{M}_t}(m|s) \mathbb{E}^{\Gamma^*, \pi^*, r^*}[U(a(h^t), a, a^{\geq t+1}, v)|h_A^t, \mathbf{M}_t, 1, m, s, a].$$

Note that the PBE assessment already specified what the agent would have done when her type is $v$ after she reported $m \in M^{+\mathbf{M}_t}$, so we only need to choose her strategy at $(h_A^t, \mathbf{M}'_t, 1)$.

For such a type $v \in V$, however, it may no longer be optimal to participate in the mechanism when the principal offers $\mathbf{M}'_t$. Thus, set $\pi_v^*(h_A^t, \mathbf{M}'_t) = 1$ only if the agent's payoff from participating is at least the payoff from not participating. Note that since we

only made worse the mechanism at $h^t$ for the agent when her type has zero probability at $h^t$, this does not affect her incentives at earlier histories. Hence, this modification does not alter the PBE, nor the payoffs or the distribution over allocations at each continuation history from the perspective of the principal and those types that have positive probability at $h^t$. It does alter the payoff and the distribution over allocations for the agent when her type has zero probability at $h^t$; however, this only happens at an event that has zero probability for her given her type. □

*Proof of Proposition A.3.* Take any $h^t$, $\mathbf{M}_t \in \operatorname{supp} \Gamma^*(h^t)$ such that there exists $s_t, s'_t \in S^{\mathbf{M}_t}$ with $\mu^*(h^t, \mathbf{M}_t, 1, s_t) = \mu^*(h^t, \mathbf{M}_t, 1, s'_t)$, where $\mu^*(h^t, \mathbf{M}_t, 1, s_t) \in \Delta(V \times H_A^{t+1}(h^t, \mathbf{M}_t, 1, s_t))$. Note that by construction, the belief does not depend on the agent's private history. In what follows, we abuse notation and denote by $\mu$ the marginal distribution on the agent's type.

The finite support assumption implies that there is $K \geqslant 1$ such that we can index the principal's posteriors at history $(h^t, \mathbf{M}_t, 1, \cdot)$ as follows $\{\mu_1, \ldots, \mu_K\}$. Partition $S^{*\mathbf{M}_t}$ as follows:

$$S^{*\mathbf{M}_t} = \bigcup_{l=1}^{K} \{s_t \in S^{*\mathbf{M}_t} : \mu^*(h^t, \mathbf{M}_t, 1, s_t) = \mu_l\} = \bigcup_{l=1}^{K} S^{*\mathbf{M}_t}(\mu_l).$$

Item 2 in Lemma A.1 implies that all the output messages in $S^{*\mathbf{M}_t}$ are generated with positive probability (from the point of view of the principal).

For each $l \in \{1, \ldots, K\}$, let $S^{\mathbf{M}_t}(\mu_l) = \{s_{t,1}^{\mu_l}, \ldots, s_{t,H_l}^{\mu_l}\}$.

Consider the following mechanism: $\mathbf{M}'_t = (\langle \beta^{\mathbf{M}'_t}, M^{\mathbf{M}'_t}, S^{\mathbf{M}'_t} \rangle, \alpha^{\mathbf{M}'_t}$, where $M^{\mathbf{M}'_t} = M^{\mathbf{M}_t}, S^{\mathbf{M}'_t} = S^{\mathbf{M}_t}$. For each $l \in \{1, \ldots, K\}$, let

$$\beta^{\mathbf{M}'_t}(s_{t,1}^{\mu_l}|m) = \sum_{h=1}^{H_l} \beta^{\mathbf{M}_t}(s|m), \quad \beta^{\mathbf{M}'_t}(s_{t,h}^{\mu_l}|m) = 0 \quad h \in \{2, \ldots, H_l\}$$

$$\alpha^{\mathbf{M}'_t}(\cdot|s_{t,1}^{\mu_l}) = \sum_{h=1}^{H_l} \frac{Pr_{\mu^*,\Gamma^*,\pi^*,r^*}(s_{t,h}^{\mu_l})}{\sum_{h'=1}^{H_l} Pr_{\mu^*,\Gamma^*,\pi^*,r^*}(s_{t,h'}^{\mu_l})} \alpha^{\mathbf{M}_t}(\cdot|s_{t,h}^{\mu_l}),$$

where for $h \in \{1, \ldots, H_l\}$,

$$Pr_{\mu*, \Gamma*, \pi*, r*}(s_{t,h}^{\mu_l}) = \sum_{(v, h_A^t)} \mu*(h^t)(v, h_A^t) \sum_{m \in M^{\mathbf{M}_t}} r_v^*(h_A^t, \mathbf{M}_t, 1)(m) \beta^{\mathbf{M}_t}(s_{t,h}^{\mu_l}|m),$$

where we are using that $\mu*(h^t)(v, h_A^t) > 0$ implies that the agent participates with probability 1.

Modify the continuation strategies as follows:

First, for those types $v \in V$ such that $\sum_{h_A^t} \mu*(h^t)(v, h_A^t) > 0$, let $\pi_v^*(h_A^t, \mathbf{M}_t') = \pi_v^*(h_A^t, \mathbf{M}_t) = 1$ and $r_v^*(h_A^t, \mathbf{M}_t', 1) = r_v^*(h_A^t, \mathbf{M}_t, 1)$. Because the original strategies do not depend on $h_A^t$ beyond $h^t$, the new strategies inherit this feature. We modify the participation and reporting strategy of the types that have zero probability at $h^t$ at the end since their strategies do not matter for the principal's incentives.

Second, for each $l \in \{1, \ldots, K\}$ and each $a \in \mathcal{A}(h^t)$ such that $\sum_{h=1}^{H_l} \alpha^{\mathbf{M}_t}(a|s_{t,h}^{\mu_l}) > 0$, partition $[0,1] = \cup_{h=0}^{H_l - 1} [\omega_h^a, \omega_{h+1}^a)$, where $\omega_0^a = 0, \omega_{H_l}^a = 1$ and for $h = 1, ..., H_l$:

$$\omega_h^a - \omega_{h-1}^a = \frac{Pr_{\mu*, \Gamma*, \pi*, r*}(s_{t,h}^{\mu_l}) \alpha^{\mathbf{M}_t}(a|s_{t,h}^{\mu_l})}{\sum_{h'=1}^{H_l} Pr_{\mu*, \Gamma*, \pi*, r*}(s_{t,h'}^{\mu_l}) \alpha^{\mathbf{M}_t}(a|s_{t,h'}^{\mu_l})}.$$

Fix any $a \in \text{supp } \alpha^{\mathbf{M}_t'}(\cdot|s_{t,\mu}^1)$. Then, for $\omega \in [\omega_{h-1}^a, \omega_h^a), h \geq 1$, let

$$(\Gamma^*, (\pi_v^*, r_v^*)_{v \in V}) \Big|_{(h^t, \mathbf{M}_t', s_{t,1}^{\mu_l}, a, \omega)} = (\Gamma^*, (\pi_v^*, r_v^*)_{v \in V}) \Big|_{(h^t, \mathbf{M}_t, s_{t,h}^{\mu_l}, a, \frac{\omega - \omega_{h-1}^a}{\omega_h^a - \omega_{h-1}^a})}.$$

That is we append to history $(h^t, \mathbf{M}_t', 1, s_{t,1}^{\mu_l}, a, \omega)$, the continuation strategy that corresponds to $(h^t, \mathbf{M}_t, 1, s_{t,h}^{\mu_l}, a, \frac{\omega - \omega_{h-1}}{\omega_h - \omega_{h-1}})$.

This clearly guarantees that the principal's payoff is the same, that he updates to $\mu_l$ after observing $s_{t,1}^{\mu_l}$ for $l \in \{1, \ldots, K\}$, and that the continuation strategies are indeed sequentially rational for the principal. Next we show that payoffs remain the same for those types of the agent to which the principal assigns positive probability, so that the above strategies are sequentially rational for them.

Fix $l \in \{1, \ldots, K\}$. For $h \in \{2, \ldots, H_l\}$, let $k_{1h}^l$ denote the following ratio:

$$k_{1h}^l = \frac{Pr_{\mu*,\Gamma*,\pi*,r*}(s_{t,h}^{\mu_l})}{Pr_{\mu*,\Gamma*,\pi*,r*}(s_{t,1}^{\mu_l})}.$$

and let $k_{11}^l = 1$. Because the principal updates to the same belief about the agent's type after each $s \in S^{\mathbf{M}_t}(\mu)$, we have that for all $v_j \in \text{supp } \mu$ and for all $h \in \{1, \ldots, |S^{\mathbf{M}_t}(\mu)|\}$,

$$\sum_{m \in M^{\mathbf{M}_t}} r_v^*(h_A^t, \mathbf{M}_t, 1)(m)\beta(s_{t,h}^{\mu_l}|m) = k_{1h}^l \left[\sum_{m \in M^{\mathbf{M}_t}} r_v^*(h_A^t, \mathbf{M}_t, 1)(m)\beta(s_{t,1}^{\mu_l}|m)\right].$$

The above expression implies that we can write

$$\alpha^{\mathbf{M}_t'}(\cdot|s_{t,1}^{\mu_l}) = \sum_{h=1}^{H_l} \frac{Pr_{\mu*,\Gamma*,\pi*,r*}(s_{t,h}^{\mu_l})}{\sum_{h'=1}^{H_l} Pr_{\mu*,\Gamma*,\pi*,r*}(s_{t,h'}^{\mu_l})}\alpha^{\mathbf{M}_t}(\cdot|s_{t,h}^{\mu_l}) = \sum_{h=1}^{H_l} \frac{k_{1h}^l}{\sum_{h'=1}^{H_l} k_{1h'}^l}\alpha^{\mathbf{M}_t}(a|s_{t,h}^{\mu_l}),$$

$$\omega_h^a - \omega_{h-1}^a = \frac{Pr_{\mu*,\Gamma*,\pi*,r*}(s_{t,h}^{\mu_l})\alpha^{\mathbf{M}_t}(a|s_{t,h}^{\mu_l})}{\sum_{h'=1}^{H_l} Pr_{\mu*,\Gamma*,\pi*,r*}(s_{t,h'}^{\mu_l})\alpha^{\mathbf{M}_t}(a|s_{t,h'}^{\mu_l})} = \frac{k_{1h}^l\alpha^{\mathbf{M}_t}(a|s_{t,h}^{\mu_l})}{\sum_{h'=1}^{H_l} k_{1h'}^l\alpha^{\mathbf{M}_t}(a|s_{t,h'}^{\mu_l})}.$$

Moreover, we also have the following:

$$\sum_{h=1}^{H_l} \frac{\sum_{m \in M^{\mathbf{M}_t}} r_v^*(h_A^t, \mathbf{M}_t, 1)(m)\beta^{\mathbf{M}_t}(s_{t,h}^{\mu_l}|m)}{\sum_{h'=1}^{H_l}(\sum_{m' \in M^{\mathbf{M}_t}} r_v^*(h_A^t, \mathbf{M}_t, 1)(m')\beta^{\mathbf{M}_t}(s_{t,h'}^{\mu_l}|m'))}\alpha^{\mathbf{M}_t}(a|s_{t,h}^{\mu_l})$$

$$= \sum_{h=1}^{H_l} \frac{k_{1h}^l}{\sum_{h'=1}^{H_l} k_{1,h'}^l}\alpha^{\mathbf{M}_t}(a|s_{t,h}^{\mu_l}) = \alpha^{\mathbf{M}_t'}(a|s_{t,1}^{\mu_l}),$$

and

$$\frac{\sum_{m \in M^{\mathbf{M}_t}} r_v^*(h_A^t, \mathbf{M}_t, 1)(m)\beta^{\mathbf{M}_t}(s_{t,h}^{\mu_l}|m)\alpha^{\mathbf{M}_t}(a|s_{t,h}^{\mu_l})}{\sum_{\tilde{h}=1}^{H_l}\sum_{\tilde{m} \in M^{\mathbf{M}_t}} r_v^*(h_A^t, \mathbf{M}_t, 1)(\tilde{m})\beta^{\mathbf{M}_t}(s_{t,\tilde{h}}^{\mu_l}|\tilde{m})\alpha^{\mathbf{M}_t}(a|s_{t,\tilde{h}}^{\mu_l})} = \frac{k_{1h}^l\alpha^{\mathbf{M}_t}(a|s_{t,h}^{\mu_l})}{\sum_{h'=1}^{H_l} k_{1h'}^{\mu_l}\alpha^{\mathbf{M}_t}(a|s_{t,h'}^{\mu_l})} = \omega_h^a - \omega_{h-1}^a.$$

Thus, we can write the agent's payoff in the original mechanism $\mathbf{M}_t$ as follows:

$$
\sum_{m\in M^{\mathbf{M}_t}} r_v^*(h_A^t,\mathbf{M}_t,1)(m)\sum_{s\in S^{\mathbf{M}_t}}\beta^{\mathbf{M}_t}(s|m)\sum_{a\in\mathcal{A}(h^t)}\alpha^{\mathbf{M}_t}(a|s)\mathbb{E}[U(a(h^t),a,a^{\geq t+1},v)|h_A^t,\mathbf{M}_t,1,m,s,a]=
$$

$$
=\sum_{m\in M^{\mathbf{M}_t}} r_v^*(h_A^t,\mathbf{M}_t,1)(m)\sum_{l=1}^{K}\sum_{h=1}^{H_l}\beta^{\mathbf{M}_t}(s_{t,h}^{\mu_l}|m)\sum_{a\in\mathcal{A}(h^t)}\alpha^{\mathbf{M}_t}(a|s_{t,h}^{\mu_l})\mathbb{E}[U(a(h^t),a,a^{\geq t+1},v)|h_A^t,\mathbf{M}_t,1,m,s_{t,h}^{\mu_l},a]
$$

$$
=\sum_{l=1}^{K}\sum_{a\in\mathcal{A}(h^t)}\left[\sum_{h=1}^{H_l}\sum_{m\in M^{\mathbf{M}_t}} r_v^*(h_A^t,\mathbf{M}_t,1)(m)\beta^{\mathbf{M}_t}(s_{t,h}^{\mu_l}|m)\alpha^{\mathbf{M}_t}(a|s_{t,h}^{\mu_l})\mathbb{E}[U(a(h^t),a,a^{\geq t+1},v)|h_A^t,\mathbf{M}_t,1,m,s_{t,h}^{\mu_l},a]\right]
$$

$$
=\sum_{l=1}^{K}\sum_{a\in\mathcal{A}(h^t)}\left(\sum_{\tilde h=1}^{H_l}\sum_{\tilde m\in M^{\mathbf{M}_t}} r_v^*(h_A^t,\mathbf{M}_t,1)(\tilde m)\beta^{\mathbf{M}_t}(s_{t,\tilde h}^{\mu_l}|\tilde m)\alpha^{\mathbf{M}_t}(a|s_{t,\tilde h}^{\mu_l})\right)\times
$$

$$
\sum_{h=1}^{H_l}\frac{\sum_{m\in M^{\mathbf{M}_t}} r_v^*(h_A^t,\mathbf{M}_t,1)(m)\beta^{\mathbf{M}_t}(s_{t,h}^{\mu_l}|m)\alpha^{\mathbf{M}_t}(a|s_{t,h}^{\mu_l})}{\sum_{\tilde h=1}^{H_l}\sum_{\tilde m\in M^{\mathbf{M}_t}} r_v^*(h_A^t,\mathbf{M}_t,1)(\tilde m)\beta^{\mathbf{M}_t}(s_{t,\tilde h}^{\mu_l}|\tilde m)\alpha^{\mathbf{M}_t}(a|s_{t,\tilde h}^{\mu_l})}\mathbb{E}[U(a(h^t),a,a^{\geq t+1},v)|h_A^t,\mathbf{M}_t,1,m,s_{t,h}^{\mu_l},a]
$$

$$
=\sum_{l=1}^{K}\left(\sum_{h'=1}^{H_l}\sum_{m'\in M^{\mathbf{M}_t}} r_v^*(h_A^t,\mathbf{M}_t,1)(m')\beta^*(s_{t,h'}^{\mu_l}|m')\right)\sum_{a\in\mathcal{A}(h^t)}\alpha^{\mathbf{M}_t'}(a|s_{t,1}^{\mu_l})\times
$$

$$
\sum_{h=1}^{H_l}\frac{\sum_{m\in M^{\mathbf{M}_t}} r_v^*(h_A^t,\mathbf{M}_t,1)(m)\beta^{\mathbf{M}_t}(s_{t,h}^{\mu_l}|m)\alpha^{\mathbf{M}_t}(a|s_{t,h}^{\mu_l})}{\sum_{\tilde h=1}^{H_l}\sum_{\tilde m\in M^{\mathbf{M}_t}} r_v^*(h_A^t,\mathbf{M}_t,1)(\tilde m)\beta^{\mathbf{M}_t}(s_{t,h}^{\mu_l}|\tilde m)\alpha^{\mathbf{M}_t}(a|s_{t,\tilde h}^{\mu_l})}\mathbb{E}[U(a(h^t),a,a^{\geq t+1},v)|h_A^t,\mathbf{M}_t,1,m,s_{t,h}^{\mu_l},a]
$$

$$
=\sum_{m\in M^{\mathbf{M}_t}} r_v^*(h_A^t,\mathbf{M}_t',1)(m)\sum_{l=1}^{K}\beta^{\mathbf{M}_t'}(s_{t,1}^{\mu_l}|m)\sum_{a\in\mathcal{A}(h^t)}\alpha^{\mathbf{M}_t'}(a|s_{t,1}^{\mu_l})\times\mathbb{E}[U(a(h^t),a,a^{\geq t+1},v)|h_A^t,\mathbf{M}_t',1,m,s_{t,1}^{\mu_l},a],
$$

where the first equality uses the labeling of the posteriors we have used throughout the proof, the second equality is obtained by changing the order of summation, the third equality is obtained by multiplying and dividing by the probability that, conditional on the belief being $\mu_l$, allocation $a$ was obtained, the fourth equality is obtained by using the definition of $\alpha^{\mathbf{M}_t'}$ and grouping the terms that represent the total probability that the output message corresponds to belief $\mu_l$, and the final equality is obtained by realizing this rewriting corresponds to the payoff the agent obtains under mechanism $\mathbf{M}_t'$.

Therefore, the agent's incentives remain the same when her type has positive probability at $h^t$.

Finally, for those $v\in V$ such that $\sum_{h_A^t\in H_A^t(h^t)}\mu^*(h^t)(v,h_A^t)=0$, choose $r_v^*(h_A^t,\mathbf{M}_t',1)$ to

solve:

$$\max_{m \in M^{\mathbf{M}'_t}} \sum_{l=1}^{K} \beta^{\mathbf{M}'_t}(s_{t,1}^{\mu_l}|m) \sum_{a \in \mathcal{A}(h^t)} \alpha^{\mathbf{M}'_t}(a|s_{t,1}^{\mu_l}) \times \mathbb{E}^{\Gamma^*,\pi^*,r^*}[U(a(h^t),a,a^{\geq t+1},v)|h_A^t,\mathbf{M}'_t,1,m,s_{t,1}^{\mu_l},a],$$

and use the payoff of this to calculate $\pi_v^*(h_A^t, \mathbf{M}'_t)$.

It is immediate that with all these modifications the assessment remains a PBE. □

## B   Proof of Theorem 3.1

*Proof.* Let $\langle \Gamma^*, (\pi_v^*, r_v^*)_{v \in V}, \mu^* \rangle$ denote a PBE that satisfies the properties of Propositions A.1-A.3. That is, the agent's strategy only depends on her payoff-relevant type, $v \in V$, the agent participates with probability 1 when her type has positive probability, and each output message corresponds to exactly one posterior belief.

Fix $t \geq 0$, a history $h^t$, and a history $h^\tau$ on the path given $h^t$ for some $\tau \geq t$. For each $\mathbf{M}_\tau \in \text{supp } \Gamma^*(h^\tau)$, define the injective mapping:

$$\sigma(h^\tau, \mathbf{M}_\tau) : S^{\mathbf{M}_\tau} \mapsto \Delta(V)$$
$$\sigma(h^\tau, \mathbf{M}_\tau)(s) = \sum_{h_A^\tau \in H_A^\tau, m \in M^{\mathbf{M}_\tau}} \mu^*(h^\tau, \mathbf{M}_\tau, 1, s_\tau)(\cdot, h_A^\tau, m), \tag{11}$$

which is well-defined by Propositions A.1 and A.3

Using this, we can define the corresponding canonical mechanism $\mathbf{M}_\tau^C$ as follows:

$$\mathbf{M}_\tau^C = \{(V, \beta^{\mathbf{M}_\tau^C}, \Delta(V)), \alpha^{\mathbf{M}_\tau^C}\}, \tag{12}$$

where[33]

$$\beta^{\mathbf{M}_\tau^C}(\mu|v) = \sum_{m \in \mathbf{M}_\tau} \beta^{\mathbf{M}_\tau}(\sigma^{-1}(h^\tau, \mathbf{M}_\tau)(\mu)|m) r_v^*(h_A^\tau, \mathbf{M}_\tau, 1)(m), \qquad (13)$$

$$\alpha^{\mathbf{M}_\tau^C}(\mu) = \alpha^{\mathbf{M}_\tau}(\sigma^{-1}(h^\tau, \mathbf{M}_\tau)(\mu)). \qquad (14)$$

Note the construction of $\beta^{\mathbf{M}_\tau^C}$ uses the fact that the agent's reporting strategy only depends on her private type and not on the payoff-irrelevant part of her private history.

Having done this transformation for $t \leqslant \tau' \leqslant \tau$, we can map any history

$$h^\tau = (h^t, \mathbf{M}_t, 1, s_t, \dots, \mathbf{M}_\tau, 1, s_\tau, a_\tau, \omega_\tau),$$

on the path of the equilibrium strategy starting from $h^t$ to:[34]

$$h^{\tau^C} = (h^t, \mathbf{M}_t^C, 1, \sigma(h^t, \mathbf{M}_t)(s_t), a_t, \dots, \mathbf{M}_\tau^C, 1, \sigma(h^\tau, \mathbf{M}_\tau)(s_\tau), a_\tau, \omega_\tau).$$

Thus, we can define the principal's strategy so that $\Gamma'(h^{\tau^C})(\mathbf{M}_\tau^C) = \Gamma^*(h^\tau)(\mathbf{M}_\tau)$.

Given $h^\tau = (h^t, \mathbf{M}_t, 1, s_t, a_t, \dots, \mathbf{M}_\tau, 1, s_\tau, a_\tau, \omega_\tau)$ and the corresponding $h^{\tau^C}$, the set of agent histories that is consistent with $h^{\tau^C}$ is given by:

$$H_A^{\tau^C}(h^{\tau^C}) =$$
$$= \left\{ \begin{array}{c} (h_A^t, \mathbf{M}_t^C, 1, v_t, \sigma(h^t, \mathbf{M}_t)(s_t), a_t, \dots, \mathbf{M}_{\tau-1}^C, 1, v_{\tau-1}, \sigma_{\tau-1}(h^{\tau-1}, \mathbf{M}_{\tau-1})(s_{\tau-1}), a_{\tau-1}, \omega_\tau) : \\ h_A^t \in H_A^t(h^t), (v_t, \dots, v_{\tau-1}) \in V^\tau \end{array} \right\}.$$

Let $\pi_v'(h_A^{\tau^C}, \mathbf{M}_\tau^C) = \pi_v^*(h_A^\tau, \mathbf{M}_\tau) = 1$ and $r_v'(h_A^{\tau^C}, \mathbf{M}_\tau^C, 1) = \delta_v$.

Let $[V \times H_A^{\tau^C}(h^{\tau^C})]^*$ denote the set of truthful histories starting from $h^t$, i.e., those that have the agent of type $v$ report $v$ throughout $t, \dots, \tau - 1$ conditional on her participating

---

[33]Suppose that $v \in V$ has probability 0 at history $h^\tau$. We can use the agent's strategy profile to construct $\beta^{\mathbf{M}_\tau^C}(\cdot|v)$. A consequence of Lemma A.1 and Proposition A.3 is that the principal assigns probability 0 to such $v$ for all $s \in S^{*\mathbf{M}_\tau}$ and hence at all $\mu \in S^{*\mathbf{M}_\tau^C}$. In other words, $\beta^{\mathbf{M}_\tau^C}(\mu|v) > 0$ implies that $\mu(v) = 0$.

[34]Note that the agent always participates on the path of the strategy.

60

in the mechanism. With this notation at hand, let:

$$\mu'(h^{\tau^C})(v, h_A^{\tau^C}) = \mu^*(h^\tau)(v)\mathbb{1}[(v, h_A^{\tau^C}) \in [V \times H_A^{\tau^C}(h^{\tau^C})]^*]$$

It remains to check that at history $h^{\tau^C}$, when the principal offers $\mathbf{M}_\tau^C \in \text{supp } \Gamma'(h^{\tau^C})$ and the output message is $\nu$, his beliefs are

$$\sum_{h_A^\tau, m \in M^{\mathbf{M}_\tau}} \mu^*(h^\tau, \mathbf{M}_\tau, 1, \sigma^{-1}(h^\tau, \mathbf{M}_\tau)(\nu))(\cdot, h_A^\tau, m) = \nu.$$

Suppose that we have shown this for $h^\tau$ and now we show it for $h^{\tau+1} = (h^\tau, \mathbf{M}_\tau^C, 1, \nu)$. Note that the probability that the agent is of type $v$ and reports $v$ to the mechanism is:

$$
\begin{aligned}
\mu'(h^{\tau^C}, \mathbf{M}_\tau^C, 1, \nu)(v, h_A^\tau, \mathbf{M}_\tau^C, 1, v, \nu) &= \frac{\mu^*(h^\tau)(v)\beta^{\mathbf{M}_\tau^C}(\nu|v)}{\sum_{\widetilde{v}} \mu^*(h^\tau)(\widetilde{v})\beta^{\mathbf{M}_\tau^C}(\nu|\widetilde{v})} \\
&= \frac{\mu^*(h^\tau)(v)\sum_{h_A^\tau \in H_A^\tau(h^\tau)}\sum_{m \in M^{\mathbf{M}_\tau}} r_v^*(h_A^\tau, \mathbf{M}_\tau, 1)(m)\beta^{\mathbf{M}_\tau}(\sigma^{-1}(h^\tau, \mathbf{M}_\tau)(\nu)|m)}{\sum_{\widetilde{v}} \mu^*(h^\tau)(\widetilde{v})\sum_{h_A'^\tau \in H_A^\tau(h^\tau)}\sum_{m' \in M^{\mathbf{M}_\tau}} r_{\widetilde{v}}^*(h_A'^\tau, \mathbf{M}_\tau, 1)(m')\beta^{\mathbf{M}_\tau}(\sigma^{-1}(h^\tau, \mathbf{M}_\tau)(\nu)|m')} \\
&= \sum_{h_A^\tau \in H_A^\tau(h^\tau), m \in M^{\mathbf{M}_\tau}} \frac{\mu^*(h^\tau)(v)r_v^*(h_A^\tau, \mathbf{M}_\tau, 1)(m)\beta^{\mathbf{M}_\tau}(\sigma^{-1}(h^\tau, \mathbf{M}_\tau)(\nu)|m)}{\sum_{\widetilde{v}} \mu^*(h^\tau)(\widetilde{v})\sum_{h_A'^\tau \in H_A^\tau(h^\tau)}\sum_{m' \in M^{\mathbf{M}_\tau}} r_{\widetilde{v}}^*(h_A'^\tau, \mathbf{M}_\tau, 1)(m')\beta^{\mathbf{M}_\tau}(\sigma^{-1}(h^\tau, \mathbf{M}_\tau)(\nu)|m')} \\
&= \sum_{m \in M^{\mathbf{M}_\tau}} \sum_{h_A^\tau} \mu^*(h^\tau, \mathbf{M}_\tau, 1, \sigma^{-1}(h^\tau, \mathbf{M}_\tau)(\nu))(v, h_A^\tau, m) = \mu^*(h^\tau, \mathbf{M}_\tau, \sigma^{-1}(h^\tau, \mathbf{M}_\tau)(\nu))(v) \\
&= \nu(v)
\end{aligned}
$$

where the first equality uses that the agent participates with probability one and reports her type truthfully, the second equality uses the definition of $\beta^{\mathbf{M}_\tau^C}$ in equation (13), the third line is a rewriting of the second by taking the summation over $h_A^\tau \in H_A^\tau, m \in M^{\mathbf{M}_\tau}$ outside, the fourth is obtained by recognizing the expression within the summation is the principal's belief at history $h^\tau$ that the agent is of type $v$ is at history $h_A^\tau$ and submitted message $m$, conditional on the output message being $\sigma^{-1}(h^\tau, \mathbf{M}_\tau)(\nu)$, and the final line uses the definition of $\sigma$ to arrive to the desired expression.

For any $h^t$, any $\tau \geqslant t$, and any $h^\tau$ that is on the path given $h^t$, and the corresponding

$h^{\tau^C}$, if $\mathbf{M}_\tau \notin \text{supp } \Gamma'(h^{\tau^C})$, then

$$\pi'_v(h_A^{\tau^C}, \mathbf{M}_\tau) = \pi_v^*(h_A^\tau, \mathbf{M}_\tau)$$
$$r'_v(h_A^{\tau^C}, \mathbf{M}_\tau, 1) = r_v^*(h_A^\tau, \mathbf{M}_\tau, 1)$$

Note that for any $m \in M^{\mathbf{M}_\tau}, s \in S^{\mathbf{M}_\tau}$, the previous transformation will take the continuation strategy that follows $(h^\tau, \mathbf{M}_\tau, s_\tau, a_\tau, \omega_{\tau+1}) = h^{\tau+1}$ to one in which the principal offers canonical mechanisms.

For any other histories, specify the strategies as in the original game.

Note that we have not modified the outcome of the game after any history $h^t$; in particular, the new strategy profile implements the path of the original profile. Moreover, the agent does not have an incentive to lie; otherwise, she would have had a deviation in the original profile. Additionally, the principal also has no deviations; otherwise, he would have had an incentive to deviate in the original profile. This completes the proof of Theorem 3.1. □

## C  Properties of the canonical game

**Proposition C.1.** *Fix a* canonical *PBE of the mechanism-selection game* $\langle \Gamma^*, (\pi_v^*, r_v^*)_{v \in V}, \mu^* \rangle$. *Then, without loss of generality, for any public history $h^t$, there exists a canonical mechanism $\mathbf{M}_t^C$ such that*

1. *$\mathbf{M}_t^C$ maximizes the principal's payoff from a deviation at $h^t$,*

2. *$\pi_v^*(h_A^t, \mathbf{M}_t^C) = 1$,*

3. *$r_v^*(h_A^t, \mathbf{M}_t^C, 1) = \delta_v$.*

*Proof.* Fix a history $h^t$ and suppose that there exists a non-canonical mechanism $\mathbf{M}_t^*$ that maximizes the principal's payoff from a deviation. Let $\pi_v^*(h_A^t, \mathbf{M}_t^*), r_v^*(h_A^t, \mathbf{M}_t^*, 1)$ denote the agent's participation and reporting strategy upon observing the principal's choice of $\mathbf{M}_t^*$.

We make three observations:

1. In a canonical PBE the continuation strategy for the agent for any $h_A^{t+1} = (h_A^t, \mathbf{M}_t^*, ...)$ does not depend on $h_A^{t+1}$. (Recall the proof of Proposition A.1 does not rely on the public history $h^t$ being on the path.)

2. Therefore, we can use the same construction as in Propositions A.2 and A.3 to transform the mechanism and the continuation strategy to guarantee that the agent participates with probability 1 after observing $\mathbf{M}_t^*$ and each output message maps exactly to one continuation belief. Denote by $\mathbf{M}_t^{**}$ the transformed mechanism.

3. Finally, we can use the transformation in Theorem 3.1 to construct from $\mathbf{M}_t^{**}$ a *canonical* mechanism, $\mathbf{M}_t^{C*}$, in which the agent reports truthfully and the recommended beliefs for the principal are indeed the beliefs obtained via Bayesian updating.

Note that $\mathbf{M}_t^{C*}$ is an available choice for the principal. It follows from the previous observations that in the original strategy profile, we can replace the best response for the agent and the continuation strategy after the principal offers $\mathbf{M}_t^{C*}$ by those obtained in transforming $\mathbf{M}_t^*$ to $\mathbf{M}_t^{C*}$. In the new strategy profile, the principal is now indifferent between deviating to $\mathbf{M}_t^*$ and to $\mathbf{M}_t^{C*}$. $\qquad\square$

**Corollary C.1.** *If $\langle \Gamma^{*C}, (\pi_v^{*C}, r_v^{*C})_{v \in V}, \mu^{*C} \rangle$ is a canonical PBE of the canonical game, then there is an equilibrium of the mechanism-selection game $\langle \Gamma^*, (\pi_v^*, r_v^*)_{v \in V}, \mu^* \rangle$ that implements the same choices by the principal and the agent on the equilibrium path.*

## D  PROOFS OF SECTION 4

### D.1  Proof of Proposition 4.1

Compare the solution to $(\mathcal{P})$ to the solution to the following program:

$$\max_{\beta, a, y, t} \sum_{i=1}^{N} \mu_{0,i} \sum_{h=1}^{M} \beta_{i,h} \sum_{q \in Q} \alpha_h(q) [w_i(q_h, y_h(q)) + t_h] \qquad (\mathcal{A})$$

$$\text{s.t.} \begin{cases} \sum \beta_{1,h} [u_{1h} - t_h] \geq 0 \\ \sum_h (\beta_{i,h} - \beta_{k,h}) [u_{i,h} - t_h] \geq 0, (\forall i)(\forall k \in \{i-1, i+1\}) \\ y_h(q_h) \in \arg\max \sum_{i=1}^{N} \mu_{0,i} \beta_{i,h} w_i(q_h, y) \end{cases},$$

where $u_{i,h}$ is shorthand for $\sum_{q \in Q} \alpha_h(q) u_i(q_h, y_h(q))$.

**Theorem D.1.** If $u_i$ satisfies Definition 4.1, then the values of $(\mathcal{P})$ and $(\mathcal{A})$ coincide.

*Proof.* We show that the solution to $(\mathcal{A})$ satisfies all the constraints of $(\mathcal{P})$. Note first that the solution to $(\mathcal{A})$ satisfies that for all $i \geqslant 2$,

$$\sum_h \beta_{i,h}[u_{i,h} - t_h] \geqslant \sum_h \beta_{i-1,h}[u_{i,h} - t_h]$$

$$\sum_h \beta_{i-1,h}[u_{i-1,h} - t_h] \geqslant \sum_h \beta_{i,h}[u_{i-1,h} - t_h],$$

so that for all $i \geqslant 2$, we have

$$\sum_h (\beta_{i,h} - \beta_{i-1,h})(u_{i,h} - u_{i-1,h}) \geqslant 0. \tag{15}$$

To show that the statement of the theorem holds, consider $i$ and $j < i - 1$. The solution to $(\mathcal{A})$ satisfies

$$\sum_h \beta_{i,h}[u_{i,h} - t_h] \geqslant \sum_h \beta_{i-1,h}[u_{i,h} - t_h]$$

$$\sum_h \beta_{i-1,h}[u_{i-1,h} - t_h] \geqslant \sum_h \beta_{i-2,h}[u_{i-1,h} - t_h]$$

$$\cdots$$

$$\sum_h \beta_{j+1,h}[u_{j+1,h} - t_h] \geqslant \sum_h \beta_{j,h}[u_{j+1,h} - t_h].$$

Adding up, we obtain

$$\sum_{k=j+1}^{i} \sum_{h=1}^{M} (\beta_{k,h} - \beta_{k-1,h}) u_{k,h} \geqslant \sum_h (\beta_{i,h} - \beta_{j,h}) t_h. \tag{16}$$

Monotonic expectational differences together with equation (15) implies the left-hand side is bounded above by

$$\sum_{k=j+1}^{i} \sum_{h=1}^{M} (\beta_{k,h} - \beta_{k-1,h}) u_{i,h} = \sum_{h=1}^{M} (\beta_{i,h} - \beta_{j,h}) u_{i,h}. \tag{17}$$

64

Equations (16) and (17) imply

$$\sum_h \beta_{i,h}[u_{i,h} - t_h] \geqslant \sum_h \beta_{j,h}[u_{j,h} - t_h].$$

Therefore, the constraint that $i$ does not report $j < i - 1$ holds.

Similarly, consider $i$ and $j > i + 1$. The solution to ($\mathcal{A}$) satisfies

$$\sum_h \beta_{i,h}[u_{i,h} - t_h] \geqslant \sum_h \beta_{i+1,h}[u_{i,h} - t_h]$$

$$\sum_h \beta_{i+1,h}[u_{i+1,h} - t_h] \geqslant \sum_h \beta_{i+2,h}[u_{i+1,h} - t_h]$$

$$\cdots$$

$$\sum_h \beta_{j-1,h}[u_{j-1,h} - t_h] \geqslant \sum_h \beta_{j,h}[u_{j-1,h} - t_h].$$

Adding up, we obtain

$$\sum_{k=i}^{j-1} \sum_{h=1}^M (\beta_{k,h} - \beta_{k+1,h}) u_{k,h} \geqslant \sum_h (\beta_{i,h} - \beta_{j,h}) t_h. \tag{18}$$

Monotonic expectational differences together with equation (15) imply that the left-hand side is bounded above by

$$\sum_{k=i}^{j-1} \sum_{h=1}^M (\beta_{k,h} - \beta_{k+1,h}) u_{i,h} = \sum_{h=1}^M (\beta_{i,h} - \beta_{j,h}) u_{i,h}. \tag{19}$$

Equation (19) follows because equation (15) implies $\sum_{h=1}^M (\beta_{k,h} - \beta_{k+1,h}) u_{k,h}$ is decreasing in $k$.

Equations (18) and (19) imply

$$\sum_h \beta_{i,h}[u_{i,h} - t_h] \geqslant \sum_h \beta_{j,h}[u_{j,h} - t_h].$$

Therefore, the incentive constraint that $i$ does not report $j$, $j > i + 1$ holds.

Finally, because we have all incentive compatibility constraints, it follows that, when

$u_i$ satisfies Definition 4.1, the participation constraints for $i \geqslant 2$ are implied by the participation constraint for $i = 1$. $\qquad\square$

**Proposition D.1.** *The participation constraint for $i = 1$ binds in the solution to* ($\mathcal{A}$).

*Proof.* Otherwise, let $\epsilon = \beta_1 \cdot (u_1 - t)$ and consider the mechanism that charges $\tilde{t}_h = t_h + \epsilon$. All incentive constraints continue to be satisfied, the participation constraint for $i = 1$ holds, and revenue increases. $\qquad\square$

### D.2 Proof of Propositions 4.2 and 4.3

We consider program ($\mathcal{A}$) but with the following modifications:

1. The participation constraint binds for $i = 1$.

2. We write everything in terms of distribution over posteriors as opposed to communication devices.

3. We replace the principal's sequential rationality constraint by the correspondence $y_\mu^*(q) \equiv \arg\max_{y \in Y(q)} \sum_{i=1}^{N} \mu_i w_i(q, y)$ and the Bayesian plausibility constraint.

Therefore, we obtain

$$\max_{\tau, a, t} \sum_{h=1}^{M} \tau(\mu_h) \sum_{i=1}^{N} \mu_{h,i} \sum_{q \in Q} \alpha_h(q)[w_i(q, y_h^*(q)) + t_h] \qquad (\mathcal{A}')$$

$$\text{s.t.} \begin{cases} \sum_{h=1}^{M} \tau(\mu_h) \frac{\mu_{h,1}}{\mu_{0,1}}[u_{1h} - t_h] = 0 \\ \sum_{h=1}^{M} \tau(\mu_h)(\frac{\mu_{h,i}}{\mu_{0,i}} - \frac{\mu_{h,k}}{\mu_{0,k}})[u_{i,h} - t_h] \geqslant 0, (\forall i)(\forall k \in \{i - 1, i + 1\}) \\ \sum_{h=1}^{M} \tau(\mu_h)\mu_{h,i} = \mu_{0,i}, i \in \{1, \ldots, N\} \end{cases},$$

where $u_{i,h}$ is shorthand for $\sum_{q \in Q} \alpha_h(q) u_i(q, y_h^*(q))$.

Now fix an allocation $\mathfrak{a} = (\alpha, t, y)$ where $\alpha : \Delta(V) \mapsto \Delta^*(Q), t : \Delta(V) \mapsto \mathbb{R}$ and $y : Q \times \Delta(V) \mapsto \Delta^*(\cup_{q,\mu} y^*(q, \mu)), y^*(q, \mu) \equiv \arg\max_{y \in Y(q)} \sum \mu_i w_i(q, y), \text{supp } y(q, \mu) \subseteq y^*(q, \mu)$.

Consider the program

$$\max_\tau \sum_{h=1}^{M} \tau(\mu_h) \sum_{i=1}^{N} \mu_{h,i} \sum_{q \in Q} \alpha_h(q)[w_i(q, y_h(q)) + t_h] \qquad (\mathcal{A}_\mathfrak{a})$$

$$\text{s.t.} \begin{cases} \sum_{h=1}^{M} \tau(\mu_h) \frac{\mu_{h,1}}{\mu_{0,1}}[u_{1h} - t_h] = 0 \\ \sum_{h=1}^{M} \tau(\mu_h)(\frac{\mu_{h,i}}{\mu_{0,i}} - \frac{\mu_{h,k}}{\mu_{0,k}})[u_{i,h} - t_h] \geqslant 0, (\forall i)(\forall k \in \{i-1, i+1\}) \\ \sum_{h=1}^{M} \tau(\mu_h)\mu_{h,i} = \mu_{0,i}, i \in \{1, \dots, N\} \end{cases}$$

Note that not all allocations $\mathfrak{a}$ can be made incentive compatible. To address this issue, let $C_0^\mathfrak{a}$ denote the policies $\tau$ that satisfy the constraints in $(\mathcal{A}_\mathfrak{a})$. Letting $f_0^\mathfrak{a}(\tau) = \sum_{h=1}^{M} \tau(\mu_h) \sum_{i=1}^{N} \mu_{h,i}[w_i(q_h, y_h(q_h)) + t_h]$, consider the modified objective function

$$f^\mathfrak{a}(\tau) = \begin{cases} f_0^\mathfrak{a}(\tau) & \text{if } \tau \in C_0^\mathfrak{a} \\ -\infty & \text{otherwise} \end{cases} .$$

In what follows, $f^\mathfrak{a}(\tau)$ is the objective function under consideration.

In Doval and Skreta (2018a), we extend the results in Le Treust and Tomala (2017) to show that given a constrained maximization problem,[35]

$$\text{cav}_{g_1,\dots,g_K} f(\mu, \gamma_1, \dots, \gamma_K) := \sup \left\{ \sum_m \lambda_m f(\mu_m) : \begin{array}{l} \sum_m \lambda_m \mu_m = \mu, \\ \sum_m \lambda_m g_l(\mu_m) \geqslant \gamma_l, l \in \{1, \dots, r\}, \\ \sum_m \lambda_m g_l(\mu_m) = \gamma_l, l \in \{r+1, \dots, K\} \end{array} \right\},$$

$$(20)$$

where $f, g_1, \dots, g_r, g_{r+1}, \dots, g_K : \Delta(V) \mapsto \mathbb{R} \cup \{-\infty\}$ is a tuple of functions defined on $\Delta(V)$, it follows that

$$\text{cav}_{g_1,\dots,g_K} f(\mu, \gamma_1, \dots, \gamma_K) = \text{cav} f^{g_1,\dots,g_K}(\mu, \gamma_1, \dots, \gamma_K)$$

---

[35]We extend the construction in their paper for completeness given that our problem includes multiple inequality constraints and equality constraints.

where $f^{g_1,\ldots,g_K} : \Delta(V) \times \mathbb{R}^K \mapsto \mathbb{R} \cup \{-\infty\}$ is such that

$$f^{g_1,\ldots,g_K}(\mu, \gamma_1, \ldots, \gamma_K) = \begin{cases} f(\mu) & \text{if } \gamma_i \leqslant g_i(\mu), i \in \{1, \ldots, r\} \wedge \gamma_i = g_i(\mu), i \in \{r+1, \ldots, K\} \\ -\infty & \text{otherwise} \end{cases}.$$
(21)

That is, the constrained *Bayesian persuasion* problem with $r$ inequality constraints and $K - r$ equality constraints in (20) can be thought of a Bayesian persuasion problem in which the objective has domain $\Delta(V) \times \mathbb{R}^K$ as defined in (21). We use this to derive properties about the number of posteriors used in the optimal solution.

Note that $(\mathcal{A}_{\mathfrak{a}})$ is a version of this program with $r = 2N - 2$ and

$$g_i(\mu) = \left[ \frac{\mu_i}{\mu_{0,i}} - \frac{\mu_{i+1}}{\mu_{0,i+1}} \right] \sum_{q \in Q} \alpha(\mu)(q)[u_i(q, y_\mu(q)) - t(\mu)], i \in \{1, \ldots N - 1\}$$

$$g_{N-2+i}(\mu) = \left[ \frac{\mu_i}{\mu_{0,i}} - \frac{\mu_{i-1}}{\mu_{0,i-1}} \right] \sum_{q \in Q} \alpha(\mu)(q)[u_i(q, y_\mu(q)) - t(\mu)], i \in \{2, \ldots N\}$$

$$g_{2N-1}(\mu) = \frac{\mu_1}{\mu_{0,1}} \sum_{q \in Q} \alpha(\mu)(q)[u_1(q, y_\mu(q)) - t(\mu)]$$

and $\gamma_i = 0$ for all $i$. We then have the following:

**Corollary D.1.** *Suppose the value of $(\mathcal{A}_{\mathfrak{a}})$ is finite. Then, $\tau$ puts positive probability in at most $3N - 1$ beliefs.*

*Proof.* This follows from Proposition **??** in Doval and Skreta (2018a) and Carathéodory's theorem (see, e.g., Rockafellar (1970)). □

Similarly, we can construct a program $(\mathcal{M}_{\mathfrak{a}})$

$$\max_\tau \sum_{h=1}^M \tau(\mu_h) \sum_{i=1}^N \mu_{h,i} \sum_{q \in Q} \alpha_h(q)[w_i(q, y_h(q)) + u_i(q, y_h(q)) - \sum_{i+1 \leqslant l} \frac{\mu_{0,l}}{\mu_{0,i}}(u_{i+1,h} - u_{i,h})]$$

$$\text{s.t.} \begin{cases} \sum_{h=1}^M \tau(\mu_h) \left[ \frac{\mu_{h,i+1}}{\mu_{0,i+1}} - \frac{\mu_{h,i}}{\mu_{0,i}} \right] (u_{i+1,h} - u_{i,h}) \geqslant 0, & i \in \{1, \ldots, N-1\} \\ \sum_{h=1}^M \tau(\mu_h) \mu_h = \mu_0 \end{cases} \quad (\mathcal{M}_{\mathfrak{a}})$$

**Corollary D.2.** *Suppose the value of $(\mathcal{M}_{\mathfrak{a}})$ is finite. Then $\tau$ puts positive probability*

*on at most $2N - 1$ beliefs.*

Another immediate corollary is that dropping constraints from a program lowers the upper bound on the number of beliefs in the support of the solution to the program:

**Corollary D.3.** *Suppose the value of $(\mathcal{M}_\mathfrak{a})$ is finite and only $M$ constraints bind. Then, $\tau$ puts positive probability on at most $N + M$ beliefs.*

*Proof.* See Corollary **??** in Doval and Skreta (2018a). $\square$

### D.3 Example 2

The next example illustrates that even if the solution to the relaxed program, $(\mathcal{R})$, satisfies the monotonicity constraints $(M)$, it may not be a solution to the original problem.

**Example 2.** Consider the sale of a durable good example in Section 2.1, but with three types $V \equiv \{v_L, v_M, v_H\}$. We provide a parametrization of the problem such that the solution to the relaxed program $(\mathcal{R})$ has the following properties:

1. In period 1, $v_H$ buys with probability 1, $v_L$ buys with probability 0, and $v_M$ buys with positive probability (but bounded away from one).

2. The allocation satisfies the monotonicity constraints, $(M)$.

3. The communication device generates two posteriors, $\mu^{HM}, \mu^{ML}$, where[36]

$$\mu^{HM}(v_H) = \frac{v_M}{v_H}, \mu^{HM}(v_L) = 0$$
$$\mu^{ML}(v_M) = \frac{v_L}{v_M}, \mu^{ML}(v_H) = 0.$$

4. However, it is not possible to find two transfers, $t(\mu^{HM}), t(\mu^{ML})$, that satisfy that (i) $v_L$'s participation constraint binds, and (ii) both $v_M$ and $v_H$'s downward-looking incentive constraints bind.

---

[36]The reader may recognize $\mu^{HM}, \mu^{LM}$ as two of the *extremal beliefs* in Bergemann et al. (2015): $\mu^{HM}$ makes the principal indifferent between setting a price of $v_M$ or $v_H$, whereas $\mu^{ML}$ makes the principal indifferent between setting a price of $v_M$ or $v_L$. Indeed, the optimal mechanism for the principal need only put weight on the extremal beliefs. Details are available from the authors upon request.

The parameters are as follows.[37] First, the prior is given by

$$\mu_0(v_H) = 0.4637, \mu_0(v_L) = 0.1194, \mu_0(v_M) = 0.4169,$$

and is chosen so that it is a convex combination of $\mu^{HM}, \mu^{ML}$. The values for the types are

$$v_H = 4.8385, v_M = 2.5528, v_L = 0.0357,$$

and are chosen so that $v_L$'s virtual valuation is negative, whereas $v_M$'s virtual valuation is positive. Also, we set $\delta = 0.95$. With these values, we have that

$$\mu^{HM}(v_H) = 0.5276$$
$$\mu^{ML}(v_M) = 0.0140.$$

and the communication device satisfies:

$$\beta(\mu^{HM}|v_H) = 1$$
$$\beta(\mu^{ML}|v_L) = 1$$
$$\beta(\mu^{ML}|v_M) = \frac{\mu^{ML}(v_M)\tau(\mu^{ML})}{\mu_0(v_M)} = \frac{0.0140 \times 0.8789}{0.4169}$$

Because $v_L$ never buys (the monopolist in period 1 recommends a price of $v_M$ is period 2), it has to be that $t(\mu^{ML}) = 0$. To determine $t(\mu^{HM})$, we have the following two equations:

$$v_H - t(\mu^{HM}) = \beta(\mu^{HM}|v_M)(v_H - t(\mu^{HM})) + \beta(\mu^{ML}|v_M)\delta \times (v_H - v_M)$$
$$\beta(\mu^{HM}|v_M)(v_M - t(\mu^{HM})) + \beta(\mu^{ML}|v_M)\delta \times (v_M - v_M) = -t(\mu^{ML})\delta(v_M - v_M) = 0.$$

The first equality implies $t(\mu^{HM}) = v_H - \delta(v_H - v_M)$, whereas the second implies $t(\mu^{HM}) = v_M$. Hence, it is not possible to find two transfers, $t(\mu^{HM}), t(\mu^{ML})$ that satisfy that the downward looking constraints bind and implements the solution to the relaxed program.

---

[37]The Matlab code, which implements the linear program used to derive the example, is available upon request.

## E  Output messages as recommendations

*Proof of Proposition 5.1.* We prove 2 implies 3. That 1 implies 2 follows from the results of the paper. Because a straightforward mechanism is a particular case of a mechanism, it follows immediately that 3 implies 1, thus completing the proof.

Thus, consider $\langle (V, \beta, \Delta(V)), \alpha \rangle$ and $y : \Delta(V) \times Q \mapsto \cup \Delta(Y(q))$ that solves $(\mathcal{R})$, or equivalently (3). The finite support assumption implies we can label $\{\mu_1, \ldots, \mu_H\}$ the posteriors that are induced with positive probability by the mechanism. Given $h \in \{1, \ldots, H\}$, let $(QY)_h = \{(q, y) : \alpha_h(q) \times y_h(q)(y) > 0\}$ denote the pairs $(q, y)$ that are implemented when the belief is $\mu_h$. Let $(QY)^* = \cup_{h \in H} (QY)_h$. Consider now the following mechanism $\langle (V, \beta', (QY)^*), \alpha' \rangle$ and the continuation strategy $y' : (QY)^* \times Q \mapsto \cup \Delta(Y(q))$ such that $\alpha'_{(q,y)}(q') = \mathbb{1}[q' = q]$ and $y'_{(q,y)}(q)(y') = \mathbb{1}[y' = y]$. Moreover, let

$$\beta'((q, y)|v_i) = \sum_{h=1}^{H} \beta(\mu_h|v_i)\alpha_h(q)y_h(q)(y).$$

Clearly, this mechanism delivers the same payoff to the principal and the agent. We now verify that it remains incentive compatible for the principal to follow the recommendations. Fix $(q, y) \in (QY)^*$. The principal's belief upon observing the output $(q, y)$ is

$$
\begin{aligned}
\mu_{(q,y)}(v_i) &= \frac{\mu_i^0 \beta'((q, y)|v_i)}{\sum_j \mu_j^0 \beta'((q, y)|v_j)} \\
&= \frac{\mu_i^0 \sum_{h=1}^{H} \beta(\mu_h|v_i)\alpha_h(q)y_h(q)(y)}{\sum_j \mu_j^0 \sum_{h=1}^{H} \beta(\mu_h|v_j)\alpha_h(q)y_h(q)(y)} \\
&= \sum_{h \in H:(q,y)\in(QY)_h} \frac{\mu_i^0 \beta(\mu_h|v_i)}{\sum_{j'} \mu_{j'}^0 \beta(\mu_h|v_{j'})} \frac{\sum_{j'} \mu_{j'}^0 \beta(\mu_h|v_{j'})}{\sum_{h=1}^{H} \sum_j \mu_j^0 \beta(\mu_h|v_j)\alpha_h(q)y_h(q)(y)} \\
&= \sum_{h \in H:(q,y)\in(QY)_h} \mu_{h,i} \frac{\tau(\mu_h)}{\sum_{h':(q,y)\in(QY)_{h'}} \tau(\mu_{h'})},
\end{aligned}
$$

where recall that $\tau(\mu_h) = \sum_{j=1}^{N} \mu_j^0 \beta(\mu_j|v_j)$. Then the payoff of the second-period principal

when he observes $(q, y)$ and chooses $y' \in Y(q)$ can be written as

$$\sum_{h \in H:(q,y) \in (QY)_h} \frac{\tau(\mu_h)}{\sum_{h':(q,y) \in (QY)_{h'}} \tau(\mu_{h'})} \sum_{i=1}^{N} \mu_{h,i} w_i(q, y') \leqslant \sum_{h \in H:(q,y) \in (QY)_h} \frac{\tau(\mu_h)}{\sum_{h':(q,y) \in (QY)_{h'}} \tau(\mu_{h'})} \sum_{i=1}^{N} \mu_{h,i} w_i(q, y),$$

where the inequality follows from knowing that $y_h(q)(y) > 0$ for all $h$ such that $(q, y) \in (QY)_h$.

Similarly, using the expression in (3), we can write the principal's payoff from the new mechanism conditional on the output being $(q, y)$ as:

$$\sum_{h \in H:(q,y) \in (QY)_h} \frac{\tau(\mu_h)}{\sum_{h':(q,y) \in (QY)_{h'}} \tau(\mu_{h'})} \sum_{i=1}^{N} \mu_{h,i} [w_i(q, y; \mu^0) + \hat{u}_i(q, y; \mu_0)],$$

and note he has no incentive to choose another $q$, because for each $h$ such that $(q, y) \in (QY)_h$, we have that $q$ is in the set of maximizers of $\sum_{i=1}^{N} \mu_{h,i}[w_i(\cdot, y_{\mu_h}; \mu^0) + \hat{u}_i(\cdot, y_{\mu_h}; \mu_0)]$. $\square$

## F    Implementation via contracts

*Proof of Proposition 5.2.*

**Necessity:** Assume there exists $t'$ such that $(\beta, q, y)$ satisfies (DIC-P). Consider $i < j$ and $\mu, \mu'$ such that $\mu(v_i)\mu'(v_j) > 0$. Then, the following must hold:

$$u_i(q(\mu), y(\mu)) - t'(\mu) \geqslant u_i(q(\mu'), y(\mu')) - t'(\mu')$$
$$u_j(q(\mu'), y(\mu')) - t'(\mu') \geqslant u_j(q(\mu), y(\mu)) - t'(\mu),$$

which implies that

$$u_j(q(\mu'), y(\mu')) - u_j(q(\mu), y(\mu)) \geqslant u_i(q(\mu'), y(\mu')) - u_i(q(\mu), y(\mu)) \tag{22}$$

That is, letting

$$D_i(\mu', \mu) = u_i(q(\mu'), y(\mu')) - u_i(q(\mu), y(\mu)), \tag{23}$$

we need that

$$D_j(\mu', \mu) \geqslant D_i(\mu', \mu), \text{ whenever } \mu'(v_j)\mu(v_i) > 0. \tag{DIC-M}$$

Note that Assumption 1 implies $D_i(\mu, \mu')$ is strictly increasing in $i$. Thus, (DIC-M) holds with strict inequality when $i < j$.

To derive the necessary conditions for the communication device, $\beta$, suppose now that $\mu(v_i)\mu'(v_i) > 0$ for $\beta(\mu|v_i)\beta(\mu'|v_i) > 0$. Because $(\beta, q, y)$ satisfies (DIC-P) for $t'$,

$$t'(\mu) - t'(\mu') = D_i(\mu, \mu').$$

Because under Assumption 1 $D_i(\mu, \mu')$ is strictly increasing in $i$, for all $j > i$, it has to be the case that $\mu'(v_j) = 0$, and for all $j < i$, it has to be the case that $\mu(v_j) = 0$. To see this, note that if $j > i$, then $D_j(\mu, \mu') > t'(\mu) - t'(\mu')$, and hence $v_j > v_i$ can never select the allocation at $\mu'$. Likewise, if $j < i$, then $D_j(\mu, \mu') < t'(\mu) - t'(\mu')$, and hence $v_j < v_i$ can never select the allocation at $\mu$.

Moreover, if there are three beliefs $\mu, \mu', \mu''$ such that $\mu(v_i)\mu'(v_i)\mu''(v_i) > 0$ such that $u_i(q(\mu), y(\mu)) \geqslant u_i(q(\mu'), y(\mu')) \geqslant u_i(q(\mu''), y(\mu''))$ and $D_i(\mu, \mu')$ and $D_i(\mu', \mu'')$ are strictly increasing in $i$, then it has to be the case that: (i) $j > i$, then $\mu'(v_j) = \mu''(v_j) = 0$, and (ii) $j < i$, then $\mu'(v_j) = \mu(v_j) = 0$. Then, $\mu'(v_i) = 1$. It follows then that there are at most three beliefs at which $v_i$ has positive probability – if we had four or more, the ones that give intermediate utility to $v_i$ assign probability one to $v_i$. Hence, they must correspond to the same belief.

Finally, suppose $\mu(v_i)\mu'(v_i) > 0$, $D_i(\mu, \mu') > 0$ and $\mu(v_j) > 0$ for $j > i$. We now show that for all $l \in \{i + 1, \ldots, j - 1\}$, it has to be the case that $\mu(v_l) > 0$. Towards a contradiction, assume there exists $v_l, l \in \{i + 1, \ldots, j - 1\}$ such that $\mu(v_l) = 0$. Because all types have positive probability, there exists $\mu' : \beta(\mu'|v_l) > 0$. Because under $t'$, $(\beta, q, y)$ satisfies (DIC-P), it follows that

$$u_l(q(\mu'), y(\mu')) - t'(\mu') \geqslant u_l(q(\mu), y(\mu)) - t'(\mu).$$

Because $\mu(v_j) > 0$, we have that

$$u_j(q(\mu), y(\mu)) - t'(\mu) \geqslant u_j(q(\mu'), y(\mu')) - t'(\mu').$$

It follows that

$$D_l(\mu', \mu) \geqslant D_j(\mu', \mu),$$

and under our assumption, this inequality is strict. Monotonic expectational differences implies that

$$D_i(\mu', \mu) \geqslant D_l(\mu', \mu) \geqslant t'(\mu') - t'(\mu).$$

The above expression contradicts that $(\beta, q, y)$ satisfies (DIC-P) under $t'$, because $v_i$ would strictly prefer to select the allocation in $\mu'$ to the allocation in $\mu$.

**Sufficiency:** Suppose $(\beta, q, y)$ satisfies the assumptions in the statement of the proposition. Then, it is possible to label the beliefs $\mu^1, \ldots, \mu^M$ so that $k < l$ implies that $\overline{v}(\mu^k) \equiv \max\{v : \mu^k(v) > 0\} \leqslant \underline{v}(\mu^l) \equiv \min\{v : \mu^l(v) > 0\}$.

Set $t'(\mu^1) = u_1(q(\mu^1), y(\mu^1))$. Note that, by definition, $v_1 = \underline{v}(\mu^1)$. For $n > 1$, define recursively

$$t'(\mu^n) = u_{\underline{v}(\mu^n)}(q(\mu^n), y(\mu^n)) - (u_{\underline{v}(\mu^n)}(q(\mu^{n-1}), y(\mu^{n-1})) - t'(\mu^{n-1})). \qquad (24)$$

We now verify that $(\beta, q, y)$ together with $t'$ satisfies (DIC-P). We first check that $v_i$ is indifferent between $\mu$ and $\mu'$ whenever $\mu(v_i)\mu'(v_i) > 0$. Monotonicity of the information structure induced by $\beta$ implies that, without loss of generality, $v_i = \underline{v}(\mu) = \overline{v}(\mu')$; moreover, $\mu = \mu^l, \mu' = \mu^k$, with $k \in \{l+1, l+2\}$ and if $k = l+2$, then $v_i = \underline{v}(\mu^{l+1}) = \overline{v}(\mu^{l+1})$.

Consider first the case in which $k = l + 1$. Then, it follows from equation (24) that

$$t'(\mu^{l+1}) - t'(\mu^l) = u_i(q(\mu^{l+1}), y(\mu^{l+1})) - u_i(q(\mu^l), y(\mu^l)),$$

so $v_i$ is indeed indifferent. Now consider the case in which $k = l + 2$, then recalling that

$\underline{v}(\mu^{l+1}) = \underline{v}(\mu^{l+2}) = v_i$, we have

$$t'(\mu^{l+1}) - t'(\mu^l) = u_i(q(\mu^{l+1}), y(\mu^{l+1})) - u_i(q(\mu^l), y(\mu^l)),$$
$$t'(\mu^{l+2}) - t'(\mu^{l+1}) = u_i(q(\mu^{l+2}), y(\mu^{l+2})) - u_i(q(\mu^{l+1}), y(\mu^{l+1})),$$

so that $v_i$ is indifferent between selecting the outcome that corresponds to either of the three beliefs, $\mu^l, \mu^{l+1}, \mu^{l+2}$. Finally, we show that when the agent is of type $v_i$, she does not want to announce any other belief $\mu^k$ such that $\mu^k(v_i) = 0$. To see this, let $\mu^l(v_i) > 0$ and consider the case in which $l < k$. Note first that

$$t'(\mu^l) = u_1(q(\mu^1), y(\mu^1)) + \sum_{n=2}^{l} (u_{\underline{v}(\mu^n)}(q(\mu^n), y(\mu^n)) - u_{\underline{v}(\mu^n)}(q(\mu^{n-1}), y(\mu^{n-1}))),$$

so that

$$t'(\mu^k) - t'(\mu^l) = \sum_{n=l+1}^{k} (u_{\underline{v}(\mu^n)}(q(\mu^n), y(\mu^n)) - u_{\underline{v}(\mu^n)}(q(\mu^{n-1}), y(\mu^{n-1})))$$
$$= \sum_{n=l+1}^{k} D_{\underline{v}(\mu^n)}(\mu^n, \mu^{n-1}).$$

Then,

$$u_i(q(\mu^l), y(\mu^l)) - u_i(q(\mu^k), y(\mu^k)) + t'(\mu^k) - t'(\mu^l) =$$
$$= u_i(q(\mu^l), y(\mu^l)) - u_i(q(\mu^k), y(\mu^k)) + \sum_{n=l+1}^{k} D_{\underline{v}(\mu^n)}(\mu^n, \mu^{n-1})$$
$$\geqslant u_i(q(\mu^l), y(\mu^l)) - u_i(q(\mu^k), y(\mu^k)) + \sum_{n=l+1}^{k} D_i(\mu^n, \mu^{n-1}) = 0,$$

where the inequality follows from (DIC-M). A similar argument shows that the same holds for $l > k$.

$\square$

REFERENCES

AKBARPOUR, M. AND S. LI (2018): "Credible mechanisms," in *Proceedings of the 2018 ACM Conference on Economics and Computation*, ACM, 371–371.

BALDER, E. J. (2001): "On ws-convergence of product measures," *Mathematics of Operations Research*, 26, 494–518.

BERGEMANN, D., B. BROOKS, AND S. MORRIS (2015): "The limits of price discrimination," *American Economic Review*, 105, 921–57.

BEST, J. AND D. QUIGLEY (2017): "Persuasion for the long run," *Working Paper*.

BESTER, H. AND R. STRAUSZ (2000): "Imperfect commitment and the revelation principle: the multi-agent case," *Economics Letters*, 69, 165–171.

——— (2001): "Contracting with imperfect commitment and the revelation principle: the single agent case," *Econometrica*, 69, 1077–1098.

——— (2007): "Contracting with imperfect commitment and noisy communication," *Journal of Economic Theory*, 136, 236–259.

BOLESLAVSKY, R. AND K. KIM (2018): "Bayesian persuasion and moral hazard," Tech. rep.

BULOW, J. I. (1982): "Durable-goods monopolists," *Journal of political Economy*, 90, 314–332.

BURGUET, R. AND J. SAKOVICS (1996): "Reserve prices without commitment," *Games and Economic Behavior*, 15, 149–164.

CAILLAUD, B. AND C. MEZZETTI (2004): "Equilibrium reserve prices in sequential ascending auctions," *Journal of Economic Theory*, 117, 78–95.

CRAWFORD, V. P. AND J. SOBEL (1982): "Strategic information transmission," *Econometrica: Journal of the Econometric Society*, 1431–1451.

DASGUPTA, P., P. HAMMOND, AND E. MASKIN (1979): "The implementation of social choice rules: Some general results on incentive compatibility," *The Review of Economic Studies*, 46, 185–216.

DEB, R. AND M. SAID (2015): "Dynamic screening with limited commitment," *Journal of Economic Theory*, 159, 891–928.

DOVAL, L. AND V. SKRETA (2018a): "Constrained information design: Toolkit," Tech. rep.

——— (2018b): "Sequentially Optimal Mechanisms: Infinite horizon," Tech. rep.

FREIXAS, X., R. GUESNERIE, AND J. TIROLE (1985): "Planning under incomplete information and the ratchet effect," *The review of economic studies*, 52, 173–191.

GEORGIADIS, G. AND B. SZENTES (2018): "Optimal Monitoring Design," Tech. rep., mimeo.

GERARDI, D. AND L. MAESTRI (2018): "Dynamic contracting with limited commitment and the ratchet effect," *Working Paper*.

GERSHKOV, A., J. K. GOEREE, A. KUSHNIR, B. MOLDOVANU, AND X. SHI (2013): "On the equivalence of Bayesian and dominant strategy implementation," *Econometrica*, 81, 197–220.

GIBBARD, A. (1973): "Manipulation of voting schemes: a general result," *Econometrica: journal of the Econometric Society*, 587–601.

GOLOSOV, M. AND L. IOVINO (2016): "Social Insurance, Information Revelation, and Lack of Commitment," *Working Paper*.

GUL, F., H. SONNENSCHEIN, AND R. WILSON (1986): "Foundations of dynamic monopoly and the coase conjecture," *Journal of Economic Theory*, 39, 155 – 190.

HART, O. D. AND J. TIROLE (1988): "Contract renegotiation and Coasian dynamics," *The Review of Economic Studies*, 55, 509–540.

HENDON, E., H. J. JACOBSEN, AND B. SLOTH (1996): "The one-shot-deviation principle for sequential rationality," *Games and Economic Behavior*, 12, 274–282.

KAMENICA, E. AND M. GENTZKOW (2011): "Bayesian persuasion," *American Economic Review*, 101, 2590–2615.

KARTIK, N., S. LEE, AND D. RAPPOPORT (2017): "Single-crossing differences on distributions," Tech. rep., Working paper.

LAFFONT, J.-J. AND J. TIROLE (1987): "Comparative statics of the optimal dynamic incentive contract," *European Economic Review*, 31, 901–926.

——— (1988): "The dynamics of incentive contracts," *Econometrica: Journal of the Econometric Society*, 1153–1175.

LE TREUST, M. AND T. TOMALA (2017): "Persuasion with limited communication capacity," *Working Paper*.

LIPNOWSKI, E. AND D. RAVID (2017): "Cheap talk with transparent motives," *Working Paper*.

LIU, Q., K. MIERENDORFF, X. SHI, AND W. ZHONG (2018): "Auctions with limited commitment," *American Economic Review, forthcoming*.

MANELLI, A. M. AND D. R. VINCENT (2010): "Bayesian and dominant-strategy implementation in the independent private-values model," *Econometrica*, 78, 1905–1938.

McADAMS, D. AND M. SCHWARZ (2007): "Credible sales mechanisms and intermediaries," *American Economic Review*, 97, 260–276.

McAFEE, R. P. AND D. VINCENT (1997): "Sequentially optimal auctions," *Games and Economic Behavior*, 18, 246–276.

MILGROM, P. AND I. SEGAL (2002): "Envelope theorems for arbitrary choice sets," *Econometrica*, 70, 583–601.

MOOKHERJEE, D. AND S. REICHELSTEIN (1992): "Dominant strategy implementation of Bayesian incentive compatible allocation rules," *Journal of Economic Theory*, 56, 378 – 399.

MYERSON, R. B. (1979): "Incentive compatibility and the bargaining problem," *Econometrica: journal of the Econometric Society*, 61–73.

——— (1982): "Optimal coordination mechanisms in generalized principal–agent problems," *Journal of Mathematical Economics*, 10, 67–81.

——— (1983): "Mechanism design by an informed principal," *Econometrica: Journal of the Econometric Society*, 1767–1797.

ROCKAFELLAR, R. T. (1970): *Convex analysis*, Princeton University Press.

SALAMANCA, A. (2016): "The value of mediated communication," *Working Paper*.

SKRETA, V. (2006): "Sequentially optimal mechanisms," *The Review of Economic Studies*, 73, 1085–1111.

——— (2015): "Optimal auction design under non-commitment," *Journal of Economic Theory*, 159, 854–890.

STOKEY, N. L. (1981): "Rational expectations and durable goods pricing," *The Bell Journal of Economics*, 112–128.

STRAUSZ, R. (2003): "Deterministic mechanisms and the revelation principle," *Economics Letters*, 79, 333–337.

STRULOVICI, B. (2017): "Contract negotiation and the Coase conjecture: A strategic foundation for renegotiation-proof contracts," *Econometrica*, 85, 585–616.

VARTIAINEN, H. (2013): "Auction design without commitment," *Journal of the European Economic Association*, 11, 316–342.

## Supplement to
## Mechanism with Limited Commitment

### I   TWO-PERIOD MODEL

Appendix I shows how to cast the model in Bester and Strausz (2001) and Bester and Strausz (2007) in the language of this paper (Section 4 does so for the case of transferable utility and preferences that satisfy increasing differences in distributions). This exercise is useful because one of the ingenuities of their model is that it can capture in the same setting mechanism design with limited commitment, but also delegation and renegotiation. Hence, understanding how the logic behind our results translates to this case as well has additional value. To facilitate the comparison between both papers, we follow their notation as much as possible.

Their environment is as follows. There is a principal and an agent. The agent has types $v \in \{v_1, ..., v_N\}$, each with probability $\mu_i^0 > 0$. The principal can commit to a mechanism $\langle (M, \beta, S), \alpha \rangle$, where $\beta : M \mapsto \Delta^*(S)$ is the communication device and the map $\alpha : S \mapsto \Delta^*(A)$ determines the allocation.[38] The output message $s$ and the allocation $a$ are publicly observed. The model captures limited commitment as follows. After observing $s$ and $a$, the principal updates his beliefs and chooses (possibly at random) $y \in Y(a)$; he cannot commit ex-ante to this choice. We endow the principal with a collection $(M_i, S_i)_{i \in \mathcal{I}}$ of input and output message sets in which each $M_i$ is finite, $|V| \leqslant |M_i|$, and $\Delta(M_i) \subseteq S_i$.[39] Moreover, we assume $(V, \Delta(V))$ is an element in that collection. Denote by $\mathcal{M}$ the set of all mechanisms with message sets $(M_i, S_i)_{i \in \mathcal{I}}$. A mechanism is *canonical* if $(V, \Delta(V))$ are its sets of input and output messages. Let $\mathcal{M}^C$ denote the set of canonical mechanisms and let $\mathbf{M}^C$ denote an element in that set.

When the agent's type is $v_i$, the allocation is $a$, and the uncommitted action is $y$, the agent obtains utility $u_i(a, y)$, while the principal obtains a payoff of $w_i(a, y)$. The above ingredients define an extensive form game as in Section 2, where timing is as follows:

- The agent observes her type, $v_i$.

- The principal offers a mechanism $\mathbf{M} = \langle (M^{\mathbf{M}}, \beta^{\mathbf{M}}, S^{\mathbf{M}}).\alpha^{\mathbf{M}} \rangle \in \mathcal{M}$.

- The agent observes the mechanism and decides whether to participate. Let $\pi_v(\mathbf{M}) \in [0, 1]$ denote the probability that type $v$ participates.[40]

  ⋆ If the agent does not participate, an allocation $a^*$ is determined. Having observed the agent's decision not to participate and $a^*$, the principal chooses $\gamma(\mathbf{M}, 0) \in \Delta^*(Y(a^*))$.

- If the agent participates, she privately chooses $m$ according to $r_v(\mathbf{M}, 1) \in \Delta(M)$.

- $s$ is drawn according to $\beta^{\mathbf{M}}(\cdot|m)$. and $a$ is drawn according to $\alpha^{\mathbf{M}}(\cdot|s)$.

- The principal observes $s$ and $a$ and chooses $\gamma(\mathbf{M}, 1, s, a) \in \Delta^*(Y(a))$.

---

[38]Bester and Strausz (2007) actually do not allow for randomized allocations. It follows from the discussion in Section 2.2 in the paper that not allowing for randomized allocations may be with loss of generality.

[39]Technically, we only need that $S_i$ contains an image of $\Delta(M_i)$.

[40]Bester and Strausz (2007) do not model participation explicitly; rather, they solve for the optimal mechanism that guarantees the agent receives at least a non-negative payoff.

We are interested in characterizing the outcomes of the game, and we use Perfect Bayesian equilibrium (PBE) as a solution concept. That is, a strategy profile $(\Gamma, \gamma, (\pi_v, r_v)_{v \in V})$ and a system of beliefs $\mu$ such that the strategies are sequentially rational given the beliefs and the beliefs are derived from Bayes' rule whenever possible. That is, letting the principal's belief about the agent's type after observing that the agent did not participate be denoted $\mu(\mathbf{M}, 0)(v_i)$, we have

$$[\sum_{v \in V}(1 - \pi_v(\mathbf{M}))]\mu(\mathbf{M}, 0)(v_i) = \mu_i^0(1 - \pi_{v_i}(\mathbf{M})).$$

Letting $\mu(\mathbf{M}, 1, s, a)(v_i)$ denote the principal's belief that the agent is of type $v_i$ when he offers mechanism $\mathbf{M}$, the agent participates, output $s$ is realized, and allocation $a$ is drawn, we have

$$[\sum_{v \in V} \sum_{m \in M^{\mathbf{M}}} \pi_v(\mathbf{M})r_v(\mathbf{M}, 1)(m)\beta^{\mathbf{M}}(s|m)]\mu(\mathbf{M}, 1, s, a)(v_i) = \mu_i^0 \pi_{v_i}(\mathbf{M}) \sum_{m \in M^{\mathbf{M}}} r_{v_i}(\mathbf{M}, 1)(m)\beta^{\mathbf{M}}(s|m).$$

To simplify notation in what follows, let $u_i(a, \gamma) = \sum_{y \in Y(a)} \gamma(y)u_i(a, y)$; define $w_i(a, \gamma)$ similarly.

We start by observing that, without loss of generality, the agent always participates in any PBE. To see this, consider a PBE assessment $\langle \Gamma, \gamma, (\pi_v, r_v)_{v \in V}), \mu \rangle$. Let $\mathbf{M} \in \text{supp } \Gamma$ be such that $\pi_v(\mathbf{M}) < 1$ for some $v \in V$. Note that the finite support assumptions on the communication device imply that there exists $s^* \in S^{\mathbf{M}}$ such that $\sum_{m \in M^{\mathbf{M}}} \beta^{\mathbf{M}}(s^*|m) = 0$. Moreover, because $|V| \leqslant |M^{\mathbf{M}}|$, label $M^{\mathbf{M}} = \{m_1, \ldots, m_N, \ldots, m_{|M^{\mathbf{M}}|}\}$. Define

$$\widetilde{\beta^{\mathbf{M}}}(s|m_i) = \sum_{m \in M^{\mathbf{M}}} \beta^{\mathbf{M}}(s|m)r_{v_i}(\mathbf{M}, 1)(m).$$

Consider then the following mechanism $\mathbf{M}' = \langle (M^{\mathbf{M}}, \beta^{\mathbf{M}'}, S^{\mathbf{M}'}), \alpha^{\mathbf{M}'} \rangle$:

$$\beta^{\mathbf{M}'}(s|m_i) = \begin{cases} \pi_{v_i}(\mathbf{M})\widetilde{\beta^{\mathbf{M}}}(s|m_i) & \text{if } s \neq s^* \\ 1 - \pi_{v_i}(\mathbf{M}) & \text{if } s = s^* \end{cases},$$

and $\alpha^{\mathbf{M}'}(s^*) = \delta_{a^*}$. Moreover, let $\gamma(\mathbf{M}', 1, s^*, a^*) = \gamma(\mathbf{M}, 0, a^*)$. Note that if we modify the strategies so that $r_{v_i}(\mathbf{M}', 1) = \delta_{m_i}$, the principal and the agent receive the same payoff, and the agent finds it optimal to participate with probability 1.

Bester and Strausz (2007) show that, without loss of generality, $M \equiv \{v_1, ..., v_N\}$ and the agent reports truthfully. Analogously, it follows that, without loss of generality, we can focus on Perfect Bayesian equilibria where the principal offers $M = V$ and the agent is truthful. To see this, consider a PBE assessment $\langle \Gamma, \gamma, (\pi_v, r_v)_{v \in V}), \mu \rangle$: Let $\mathbf{M} \in \text{supp } \Gamma$ and $r_{v_i}(\mathbf{M}, 1) \in \Delta(M^{\mathbf{M}})$ denote the reporting strategy of the agent when her type is $v_i$ and the principal offers $\mathbf{M}$. Then, her payoff can be written as

$$\sum_{m \in M^{\mathbf{M}}} r_{v_i}(\mathbf{M}, 1)(m) \sum_{s \in S^{\mathbf{M}}} \beta^{\mathbf{M}}(s|m) \sum_{a \in A} \alpha^{\mathbf{M}}(a|s) u_i(a, \gamma(\mathbf{M}, s, a))$$

$$= \sum_{s \in S^{\mathbf{M}}} \left( \sum_{m \in M^{\mathbf{M}}} r_{v_i}(\mathbf{M}, 1)(m) \beta^{\mathbf{M}}(s|m) \right) \sum_{a \in A} \alpha^{\mathbf{M}}(a|s) u_i(a, \gamma(\mathbf{M}, s, a))$$

$$= \sum_{s \in S^{\mathbf{M}}} \beta^*(s|v_i) \sum_{a \in A} \alpha^{\mathbf{M}}(a|s) u_i(a, \gamma(\mathbf{M}, s, a)).$$

Therefore, by selecting $\langle (V, \beta^*, S), \alpha \rangle$, the principal can implement the same outcome and truthtelling is the agent's best response. We focus hereafter on equilibria of the game in which the principal chooses mechanisms with $M = V$ and the agent participates with probability one and truthfully reports her type on the equilibrium path.

We now argue, that without loss of generality, if $s_h \neq s_j$, then $\mu^h \neq \mu^j$. Hence, output messages can be taken to be the principal's posterior beliefs about the agent's type. Consider, for example, a PBE and let $\mathbf{M} \in \text{supp } \Gamma$. Suppose $s_h \neq s_j$ exist, but $\mu \equiv \mu(\mathbf{M}, 1, s_h, \cdot) = \mu(\mathbf{M}, 1, s_j, \cdot)$.[41] Let $S_\mu = \{s_1^\mu, \ldots, s_{H_\mu}^\mu\}$ denote the set of output messages that lead to belief $\mu$. One can alternatively define the following communication device and allocation rule:

$$\beta'(s_1^\mu|v) = \sum_{s \in S_\mu} \beta^{\mathbf{M}}(s|v), \beta'(s_h^\mu|v) = 0, h \in \{2, \ldots, H_\mu\}$$

$$\alpha'(a|s_1^\mu) = \sum_{s \in S_\mu} \frac{\sum_{i=1}^{N} \mu_i^0 \beta^{\mathbf{M}}(s|v_i)}{\sum_{s' \in S_{\mu^h}} \sum_{i=1}^{N} \mu_i^0 \beta^{\mathbf{M}}(s'|v_i)} \alpha^{\mathbf{M}}(a|s),$$

$$\gamma'(\mathbf{M}', 1, a, s_1^\mu)(y) = \sum_{s \in S_{\mu^h}} \frac{\sum_{i=1}^{N} \mu_i^0 \beta^{\mathbf{M}}(s|v_i) \alpha^{\mathbf{M}}(a|s)}{\sum_{s' \in S_{\mu^h}} \sum_{i=1}^{N} \mu_i^0 \beta^{\mathbf{M}}(s'|v_i) \alpha^{\mathbf{M}}(a|s')} \gamma(\mathbf{M}, 1, s, a)(y),$$

---

[41] Recall that updating -conditional on participation-does not depend on the allocation $a$.

where $\beta', \alpha', \gamma'$ coincide with the original mechanism and continuation strategy for the remaining output messages, and we let $\mathbf{M}'$ in the definition of the new strategy for the principal be the mechanism defined by the above communication device and allocation rule. Clearly this does not change the principal's payoff, and when he observes $s_1^\mu$, his beliefs are updated to $\mu$, such that $\gamma'$ is optimal.

We now verify that the agent's incentive to tell the truth remains the same. Toward this end, for $i \in \{2, \ldots, H_\mu\}$, let $k_i$ denote the following ratio:

$$k_{1i} = \frac{\sum_{j=1}^{N} \beta(s_i|v_j)\mu_j^0}{\sum_{j=1}^{N} \beta(s_1|v_j)\mu_j^0},$$

and let $k_{11} = 1$. Because all output messages in $S_\mu$ lead to a belief of $\mu$, we have that for all $v_j$ in the support of $\mu$,

$$\beta(s_i|v_j) = k_{1i}\beta(s_1|v_j).$$

Note that $k_1.$ is independent of $v_j$. Using this we can write

$$\beta'(s_1^\mu|v) = \beta^{\mathbf{M}}(s_1^\mu|v) \sum_{i=1}^{H_\mu} k_{1i}$$

$$\alpha'(a|s_1^\mu) = \sum_{h=1}^{H_\mu} \frac{k_{1h}}{\sum_{h'=1}^{H_\mu} k_{1h'}} \alpha^{\mathbf{M}}(a|s_h),$$

$$\gamma'(\mathbf{M}', 1, s_1^\mu, a)(y) = \sum_{h=1}^{H_\mu} \frac{k_{1h}\alpha^{\mathbf{M}}(a|s_h)}{\sum_{h'=1}^{H_\mu} k_{1h'}\alpha^{\mathbf{M}}(a|s_{h'})} \gamma(\mathbf{M}, 1, s_h, a)(y).$$

Thus, we can write the agent's utility when her type is $v_j \in \operatorname{supp} \mu$ under the new mechanism as follows

$$\beta(s_1^\mu|v_j)[\sum_{a \in A} \alpha'(a|s_1^\mu) \sum_y \gamma'(\mathbf{M}', 1, s_1^\mu, a)(y)u_j(a, y)]$$

$$+ \sum_{s \notin S_\mu} \beta^{\mathbf{M}}(s|v_j) \sum_{a \in A} \alpha^{\mathbf{M}}(a|s) \sum_{y \in Y} \gamma(\mathbf{M}, 1, s, a)(y)u_j(a, y),$$

where we can expand the term in brackets to obtain

$$\beta(s_1^\mu|v_j)[\sum_{a\in A}\alpha'(a|s_1^\mu)\sum_y\left(\sum_{h=1}^{H_\mu}\frac{k_{1h}\alpha^{\mathbf{M}}(a|s_h)}{\sum_{h'=1}^{H_\mu}k_{1h'}\alpha^{\mathbf{M}}(a|s_{h'})}\gamma(\mathbf{M},1,s_h,a)(y)\right)u_j(a,y)]=$$

$$\beta(s_1^\mu|v_j)[\sum_{a\in A}\alpha'(a|s_1^\mu)\left(\sum_{h=1}^{H_\mu}\frac{k_{1h}\alpha^{\mathbf{M}}(a|s_h)}{\sum_{h'=1}^{H_\mu}k_{1h'}\alpha^{\mathbf{M}}(a|s_{h'})}\right)\left(\sum_{y\in Y}\gamma(\mathbf{M},1,s_h,a)(y)u_j(a,y)\right)]=$$

$$\beta(s_1^\mu|v_j)[\sum_{a\in A}\left(\sum_{h=1}^{H_\mu}\frac{k_{1h}}{\sum_{h'=1}^{H_\mu}k_{1h'}}\alpha^{\mathbf{M}}(a|s_h)\right)\left(\sum_{h=1}^{H_\mu}\frac{k_{1h}\alpha^{\mathbf{M}}(a|s_h)}{\sum_{h'=1}^{H_\mu}k_{1h'}\alpha^{\mathbf{M}}(a|s_{h'})}\right)\left(\sum_{y\in Y}\gamma(\mathbf{M},1,s_h,a)(y)u_j(a,y)\right)]$$

$$=\beta(s_1^\mu|v_j)[\sum_{a\in A}\left(\sum_{h=1}^{H_\mu}\frac{k_{1h}}{\sum_{h'=1}^{H_\mu}k_{1h'}}\alpha^{\mathbf{M}}(a|s_h)\right)\left(\sum_{y\in Y}\gamma(\mathbf{M},1,s_h,a)(y)u_j(a,y)\right)]$$

$$=\sum_{h=1}^{H_\mu}\beta(s_h|v_j)\sum_a\alpha^{\mathbf{M}}(s_h|v_j)\sum_y\gamma(\mathbf{M},1,s_h,a)(y)u_j(a,y),$$

so that the agent's utility remains the same. Similar steps also show that truthtelling is preserved. Therefore, we can take $S=\Delta(V)$, $\beta:V\mapsto\Delta^*(\Delta(V))$, $\alpha:\Delta(V)\mapsto\Delta^*(A)$.

Finally, arguments similar to those in the proof of Proposition 3.1 in the main text imply the principal's search for an optimal mechanism can be constrained to the class of canonical mechanisms.

Thus, we can write the principal's problem as follows. Following the notation in Bester and Strausz (2007), given a canonical mechanism $\mathbf{M}^C$, label $\{\mu_1,...,\mu_H\}\subseteq\Delta(V)$ the output messages in the support of the communication device. Let $\beta_{i,h}\equiv\beta(\mu_h|v_i)$. Given $\{\mu_1,...,\mu_H\}$, we can write the principal's problem as follows:

$$\max_{\beta,a}\sum_{i,h}\mu_i^0\beta_{i,h}\sum_{a\in A}\alpha(a|s_h)w_i(a,\gamma_h(a)) \tag{25}$$

$$s.t.\begin{cases}(\forall i,i'\in\{1,\dots,N\})\sum_h\beta_{i,h}\sum_{a\in A}\alpha(a|s_h)u_i(a,\gamma_h(a))\geqslant\sum_h\beta_{i',h}\sum_{a\in A}\alpha(a|s_h)u_i(a,\gamma_h(a))\\(\forall i\in\{1,\dots,N\})\sum_h\beta_{i,h}\sum_{a\in A}\alpha(a|s_h)u_i(a,\gamma_h(a))\geqslant 0\\\text{supp }\gamma_h(a)\subseteq\arg\max_{y\in Y(a)}\sum_i\mu_i^h w_i(a,y)\\(\sum_j\beta_{j,h}\mu_j^0)\mu_i^h=\beta_{i,h}\mu_i^0\end{cases}.$$

Rewrite the objective function in (25) as follows:

$$\sum_{\mu\in\Delta(V)}\tau(\mu)\sum_{a\in A}\alpha(a|\mu)\sum_i^N\mu_iw_i(a,\gamma_\mu(a)),$$

where

$$\tau(\mu)=\sum_{v\in V}\mu^0(v)\beta(\mu|v).$$

Moreover, we can rewrite the incentive compatibility and participation constraints as follows:

$$\sum_{\mu\in\Delta(V)}\frac{\tau(\mu)\mu_i}{\mu_i^0}\left[\sum_{a\in A}\alpha(a|\mu)u_i(a,\gamma_\mu(a))\right]\geq\sum_{\mu\in\Delta(V)}\frac{\tau(\mu)\mu_{i'}}{\mu_{i'}^0}\left[\sum_{a\in A}\alpha(a|\mu)u_i(a,\gamma_\mu(a))\right]\quad(26)$$

$$\sum_{\mu\in\Delta(V)}\frac{\tau(\mu)\mu_i}{\mu_i^0}\left[\sum_{a\in A}\alpha(a|\mu)u_i(a,\gamma_\mu(a))\right]\geq0.\quad(27)$$

Then, letting $w_i(\alpha,\gamma_\mu)=\sum_{a\in A}\alpha(a)w_i(a,\gamma_\mu(a))$, we can write (25) as follows:

$$\max_{\tau,\alpha}\mathbb{E}_\tau\mathbb{E}_\mu[w.(\alpha(\mu),\gamma_\mu]\quad(28)$$

$$\text{s.t.}\begin{cases}\mathbb{E}_\tau\mu=\mu^0\\ \text{Equations (26)-(27)}\\ (\forall a\in A)\text{supp }\gamma_\mu(a)\subseteq\arg\max_{y\in Y(a)}\sum_{i=1}^N\mu_iw_i(a,y)\end{cases}.$$

We have the following:

**Proposition I.1.** *The following are equivalent:*

1. *There exists a mechanism $\langle(V,\beta,S),\alpha\rangle$ that solves (25).*

2. *There exists a canonical mechanism $\langle(V,\beta',\Delta(V)),\alpha'\rangle$ that solves (25).*

3. *There exists a Bayes' plausible distribution over posteriors and an allocation rule $\alpha:\Delta(V)\mapsto\Delta^*(A)$ that solves (28).*

The result follows immediately from the previous discussion. Proposition I.1 does not allow us to interpret beliefs merely as recommendations, because $\alpha(\mu)$ need not

maximize $\mathbb{E}_\mu w(\tilde{\alpha}, \gamma_\mu)$. As we showed in Section 4, when there is transferable utility and the agent's utility function satisfies single - crossing, this property holds for the solution to the relaxed program: The first-period principal chooses $\alpha(\mu)$ to maximize the expected virtual surplus, where expectations are taken with respect to $\mu$, whereas the second-period principal chooses $\gamma_\mu(a)$ to maximize $\sum \mu_i w_i(a, y)$. Note, however, that for each realized allocation $a$, we can think of the first period principal as sending recommendations to the second period principal that will be obeyed.

As in Section 4, if we assume $u_i$ satisfies monotonic expectational differences (see Kartik et al. (2017) or Definition 4.1 in the main text), we can reduce the principal's problem as follows:

**Proposition I.2.** *Suppose the family $(u_i)_{i=1}^N$ satisfies monotonic expectational differences. Then, to find a solution to the principal's problem, it suffices to check that*

1. *The local incentive compatibility constraint holds (for all $v_i$, the agent does not report $v_{i-1}$ or $v_{i+1}$ when her type is $v_i$).*

2. *The participation constraint holds when the agent's type is $v_1$.*

Using these conditions, and the results in Doval and Skreta (2018a), we can derive the analogous result to Propositions 4.2-4.3 in the main text:[42]

**Proposition I.3.** *Suppose the family $(u_i)_{i=1}^N$ satisfies monotonic expectational differences. Then the optimal mechanism for the principal uses at most $3N - 1$ posteriors.*

## II   GENERAL COMMUNICATION DEVICES

As discussed in footnote 2 in the main text, we could have considered the following more general version of a mechanism:

$$\mathbf{M}_t = (M^{\mathbf{M}_t}, \beta^{\mathbf{M}_t}, S^{\mathbf{M}_t}) \tag{29}$$

where $\beta^{\mathbf{M}_t} : M^{\mathbf{M}_t} \mapsto \Delta^*(S^{\mathbf{M}_t} \times A)$. That is, associated to each input message there is a joint distribution over output messages and allocations.

It is immediate that for any $\mathbf{M}'_t$ as defined in the main text, we can let $\tilde{\beta}(s, a|m) =$

---

[42]See also the discussion following Proposition 4.3.

$\beta^{\mathbf{M}'_t}(s|m)\alpha^{\mathbf{M}'_t}(a|s)$ and obtain a mechanism as in equation (29). This appendix shows that, conditional on showing that the canonical set of output messages is $\Delta(V)$, then the formulation in the main text is equivalent to that in equation (29). The formulation of a mechanism in the main text has the advantage that it highlights the role of the communication device separately from that of the allocation rule. Since this is without loss of generality, we favored the definition in the main text for "pedagogical" purposes.

We proceed as follows. Section II.1 shows that Proposition A.3, which shows that without loss of generality there is a one-to-one mapping between output messages and posterior beliefs, extends to the case in which mechanisms are defined as in equation (29).[43] Section II.2 then shows that any mechanism as in equation (29) can be written as in the main text once we know that $S^{\mathbf{M}_t} \simeq \Delta(V)$.

## II.1  Proof of Proposition A.3 for mechanisms as in equation (29)

Following the notation in the main text, given a mechanism $\mathbf{M}_t$, define

$$(S \otimes A)^{*\mathbf{M}_t} = \{(s,a) \in S^{\mathbf{M}_t} \times A : (\exists m \in M^{\mathbf{M}_t}) : \beta^{\mathbf{M}_t}(s,a|m) > 0\},$$

to be the set of pairs $(s,a)$ that are possible under mechanism $\mathbf{M}_t$. As in the main text, we remove from the tree all those public histories that are consistent with mechanism $\mathbf{M}_t$ and $(s,a) \in S^{\mathbf{M}_t} \times A \backslash (S \otimes A)^{*\mathbf{M}_t}$.

Similarly, define $S^{*\mathbf{M}_t} = \{s \in S^{\mathbf{M}_t} : (\exists a \in A) : (s,a) \in (S \otimes A)^{*\mathbf{M}_t}\}$. Consider then a Perfect Bayesian equilibrium assessment, $\langle \Gamma^*, (\pi^*_v, r^*_v)_{v \in V}, \mu^* \rangle$, that satisfies the following properties, which follow from Propositions A.1 and A.2 and Lemma A.1:

1. The agent's equilibrium strategy only depends on her type $v \in V$ and the public history,

2. For all $t$ and public histories $h^t$, for all $v \in V$, $h^t_A, h^t_A{}' \in H^t_A(h^t)$ such that $\mu^*(h^t)(v, h^t_A), \mu^*(h^t)(v, h^t_A{}'$ 0, then $\mu^*(h^t)(v, h^t_A{}') = \mu^*(h^t)(v, h^t_A)$,

3. For all $t, h^t$, for all $v \in V$ such that $\sum_{h^t_A \in H^t_A(h^t)} \mu^*(h^t)(v, h^t_A) > 0$, then $\pi^*_v(h^t_A, \mathbf{M}_t) = 1$ for all $\mathbf{M}_t \in \text{supp } \Gamma^*(h^t)$,

---

[43]It is immediate that the proofs of Propositions A.1 and A.2 and Lemma A.1 do not depend on how we defined the mechanism.

4. For all $t, h^t$, for all $\mathbf{M}_t \in \text{supp } \Gamma^*(h^t)$, if $(s, a) \in (S \otimes A)^{*\mathbf{M}_t}$, then

$$\sum_{(v, h_A^t) \in V \times H_A^t(h^t), m \in M^{\mathbf{M}_t}} \mu^*(h^t)(v, h_A^t) r_v^*(h_A^t, \mathbf{M}_t, 1)(m) \beta^{\mathbf{M}_t}(s, a|m) > 0.$$

For such an assessment, we now show that, without loss of generality, there is a one-to-one mapping between output messages $s \in S^{*\mathbf{M}_t}$ and continuation beliefs. That is, for every $t$, public history $h^t$, $\mathbf{M}_t \in \text{supp } \Gamma^*(h^t)$, if $s_t, s_t' \in S^{*\mathbf{M}}_t$ is such that $s_t \neq s_t'$, then $\mu^*(h^t, \mathbf{M}_t, 1, s_t, a) \neq \mu^*(h^t, \mathbf{M}_t, 1, s_t', a')$ for any $a, a'$ such that $(s_t, a), (s_t', a') \in (S \otimes A)^{*\mathbf{M}_t}$.

*Proof.* Fix a history $h^t$ and $\mathbf{M}_t \in \text{supp } \Gamma^*(h^t)$. The finiteness of $M^{\mathbf{M}_t}$ and the finite support assumption on $\beta^{\mathbf{M}_t}$ implies that there exists $1 \leqslant K \leqslant |(S \otimes A)^{*\mathbf{M}_t}|$ such that the principal's prior at $h^t$ splits into $H$ posteriors, $\{\mu_1, \dots, \mu_L\}$, after observing $(s, a) \in (S \otimes A)^{*\mathbf{M}_t}$. Hence, we can write,

$$(S \otimes A)^{*\mathbf{M}_t} =$$
$$= \bigcup_{k=1}^{K} ((S \otimes A)^{*\mathbf{M}_t})_k = \bigcup_{k=1}^{K} \{(s, a) \in (S \otimes A)^{*\mathbf{M}_t} : \mu^*(h^t, \mathbf{M}_t, 1, s, a)(\cdot) = \mu_k\}.$$

Let $\{s_1^*, \dots, s_K^*\}$ denote $K$ elements of $S^{\mathbf{M}_t}$. Define $\beta^{\mathbf{M}_t'} : M^{\mathbf{M}_t} \mapsto \Delta^*(S^{\mathbf{M}_t} \times A)$ as follows

$$\beta^{\mathbf{M}_t'}(s_k^*, a|m) = \sum_{(s,a) \in ((S \otimes A)^{*\mathbf{M}_t})_k} \beta^{\mathbf{M}_t}(s, a|m).$$

Define $\mathbf{M}_t' = (M^{\mathbf{M}_t}, \beta^{\mathbf{M}_t'}, S^{\mathbf{M}_t})$. At history $h^t$, let the principal offer mechanism $\mathbf{M}_t'$ instead of $\mathbf{M}_t$. In a slight abuse of notation, let the agent's best response be determined as follows. For all $v \in V$ and $h_A^t \in H_A^t(h^t)$, let $\pi_v^*(h_A^t, \mathbf{M}_t') = \pi_v^*(h_A^t, \mathbf{M}_t)$ and $r_v^*(h_A^t, \mathbf{M}_t', 1) = r_v^*(h_A^t, \mathbf{M}_t, 1)$.

Note that (i) $\beta(s_k^*, a|m) > 0 \Leftrightarrow (\exists s \in S^{\mathbf{M}_t}) : (s, a) \in ((S \otimes A)^{*\mathbf{M}_t})_k$, and (ii) when the principal observes $(s_k^*, a)$ for any $a \in A$ such that $\sum_{m \in M^{\mathbf{M}_t}} \beta(s_k^*, a|m) > 0$, his

88

beliefs update to $\mu_k$. To see that (ii) holds, note that

$$\sum_{h_A^t \in H_A^t(h^t), m \in M^{\mathbf{M}_t}} \mu^*(h^t, \mathbf{M}_t', 1, s_k^*, a)(v, h_A^t, m) =$$

$$= \sum_{h_A^t \in H_A^t(h^t), m \in M^{\mathbf{M}_t}} \frac{\mu^*(h^t)(v, h_A^t) r_v^*(h_A^t, \mathbf{M}_t', 1)(m) \beta^{\mathbf{M}_t'}(s_k^*, a|m)}{\sum_{(v', h_A^{t\,'}), m' \in M^{\mathbf{M}_t}} \mu^*(h^t)(v', h_A^{t\,'}) r_{v'}^*(h_A^{t\,'}, \mathbf{M}_t', 1)(m') \beta^{\mathbf{M}_t'}(s_k^*, a|m')}$$

$$= \sum_{(s,a) \in ((S \otimes A)^{*\mathbf{M}_t})_h} \left\{ \frac{\sum_{h_A^t \in H_A^t(h^t), m \in M^{\mathbf{M}_t}} r_v^*(h_A^t, \mathbf{M}_t, 1)(m) \beta^{\mathbf{M}_t}(s, a|m)}{\sum_{(\widetilde{v}, \widetilde{h_A^t}) \widetilde{m} \in M^{\mathbf{M}_t}} \mu^*(h^t, \mathbf{M}_t, 1, s, a) r_{\widetilde{v}}(\widetilde{h_A^t}, \mathbf{M}_t, 1)(\widetilde{m}) \beta^{\mathbf{M}_t}(s, a|\widetilde{m})} \times \right.$$

$$\left. \frac{\sum_{(\widetilde{v}, \widetilde{h_A^t}) \widetilde{m} \in M^{\mathbf{M}_t}} \mu^*(h^t, \mathbf{M}_t, 1, s, a) r_{\widetilde{v}}(\widetilde{h_A^t}, \mathbf{M}_t, 1)(\widetilde{m}) \beta^{\mathbf{M}_t}(s, a|\widetilde{m})}{\sum_{(v', h_A^{t\,'}), m' \in M^{\mathbf{M}_t}} r_{v'}^*(h_A^{t\,'}, \mathbf{M}_t', 1)(m') \beta^{\mathbf{M}_t'}(s_k^*, a|m')} \right\}$$

$$= \sum_{(s,a) \in ((S \otimes A)^{*\mathbf{M}_t})_h} \mu_k(v) \frac{\sum_{(\widetilde{v}, \widetilde{h_A^t}) \widetilde{m} \in M^{\mathbf{M}_t}} \mu^*(h^t, \mathbf{M}_t, 1, s, a) r_{\widetilde{v}}(\widetilde{h_A^t}, \mathbf{M}_t, 1)(\widetilde{m}) \beta^{\mathbf{M}_t}(s, a|\widetilde{m})}{\sum_{(v', h_A^{t\,'}), m' \in M^{\mathbf{M}_t}} r_{v'}^*(h_A^{t\,'}, \mathbf{M}_t', 1)(m') \beta^{\mathbf{M}_t'}(s_k^*, a|m')}$$

$$= \mu_k(v).$$

We now modify the continuation strategies. For each $1 \leqslant k \leqslant K$ and $a \in A$ such that $\sum_{m \in M^{\mathbf{M}_t}} \beta^{\mathbf{M}_t'}(s_k^*, a|m) > 0$, label $\{s_1^k, \ldots, s_{H_{k,a}}^k\}$ the output messages $s$ such that $(s, a) \in ((S \otimes A)^{*\mathbf{M}_t})_k$. Partition $[0, 1] = \cup_{h=0}^{H_{k,a}-1}[\omega_h^a, \omega_{h+1}^a)$, where $\omega_0^a = 0 = 1 - \omega_{H_{k,a}}$, and for $h = 1, \ldots, H_{k,a} - 1$

$$\omega_h^a - \omega_{h-1}^a = \frac{\sum_{(v, h_A^t) \in V \times H_A^t(h^t), m \in M^{\mathbf{M}_t}} \mu^*(h^t)(v, h_A^t) r_v^*(h_A^t, \mathbf{M}_t, m) \beta^{\mathbf{M}_t}(s_h^k, a|m)}{\sum_{h'=1}^{H_{k,a}} \sum_{(v', h_A^{t\,'}) \in V \times H_A^t(h^t), m' \in M^{\mathbf{M}_t}} \mu^*(h^t)(v', h_A^{t\,'}) r_{v'}^*(h_A^{t\,'}, \mathbf{M}_t, m') \beta^{\mathbf{M}_t}(s_{h'}^k, a|m')}$$

Then, modify the continuation strategies so that

$$(\Gamma^*, (\pi_v^*, r_v^*)_{v \in V})|_{(h^t, \mathbf{M}_t', 1, s_k^*, a, \omega)} = (\Gamma^*, (\pi_v^*, r_v^*)_{v \in V})|_{(h^t, \mathbf{M}_t, 1, s_h^k, a, \frac{\omega - \omega_{h-1}^a}{\omega_h^a - \omega_{h-1}^a})}.$$

Similar steps as in the proof of Proposition A.3 show that the agent's payoff has not changed and hence her the specified strategies are still a best response. $\square$

*II.2 Equivalence*

Now consider an equilibrium assessment of the mechanism-selection game in which the principal uses canonical mechanisms in each period $t$ and history $h^t$, that is

$$\mathbf{M}_t^C = (V, \beta^{\mathbf{M}_\tau^C}, \Delta(V)), \tag{30}$$

where $\beta^{\mathbf{M}_\tau^C} : V \mapsto \Delta^*(\Delta(V) \times A)$. Thus, if $\mathbf{M}_t^C \in \text{supp } \Gamma^*(h^t)$, then

$$\mu^*(h^t, \mathbf{M}_t, 1, \nu, a)(v, \cdot) = \frac{\mu^*(h^t)(v)\beta^{\mathbf{M}_\tau^C}(\nu, a|v)}{\sum_{v' \in V} \mu^*(h^t)(v')\beta^{\mathbf{M}_t^C}(\nu, a|v')} = \mu(v),$$

whenever $\sum_{v' \in V} \mu^*(h^t)(v')\beta^{\mathbf{M}_t^C}(\nu, a|v') > 0$. Fix $v \in V$ such that $\mu(h^t)(v) > 0$ and $\nu(v) > 0$. The finite support assumption allows us to label the set $\{a \in A : \beta^{\mathbf{M}_t^C}(\nu, a|v) > 0\}$ as $\{a_1, \ldots, a_{N^\nu}\}$ for some $N^\nu \in \mathbb{N}^\nu$.[44]

Define

$$k_i^\nu = \sum_{v \in V, h_A^t \in H_A^t(h^t)} \mu^*(h^t)(v, h_A^t)\beta^{\mathbf{M}_t^C}(\nu, a_i|v).$$

Note that Bayesian updating implies that for all $v \in V$ such that $\nu(v) > 0$, the following holds

$$\frac{k_i^\nu}{k_1^\nu} = \frac{\beta(\nu, a_i|v)}{\beta(\nu, a_1|v)}. \tag{31}$$

Define the following mechanism

$$\mathbf{M}_t'^C = (\langle V, \beta^{\mathbf{M}_t'^C}, \Delta(V)\rangle, \alpha^{\mathbf{M}_t'^C}), \tag{32}$$

---

[44]Note that as long as $v$ satisfies the conditions, the set does not depend on the selected type

where

$$\alpha^{\mathbf{M}_t'^C}(a|\mu) = \frac{\sum_{v \in V, h_A^t \in H_A^t(h^t)} \mu^*(h^t)(v, h_A^t) \beta^{\mathbf{M}_t^C}(\mu, a|v)}{\sum_{v' \in V, h_A^t \in H_A^t(h^t)} \mu^*(h^t)(v', h_A^t) \sum_{a' \in A} \beta^{\mathbf{M}_t^C}(\mu, a'|v')}$$

$$\beta^{\mathbf{M}_t'^C}(\mu|v) = \sum_{a \in A} \beta^{\mathbf{M}_t^C}(\mu, a|v).$$

Fix $\nu$ such that $\beta^{\mathbf{M}_t^C}(\nu, a|v) > 0$ for some $v \in \text{supp } \sum_{h_A^t} \mu^*(h^t)(\cdot, h_A^t)$. Note that

$$\begin{aligned}
\alpha^{\mathbf{M}_t'^C}(a_i|\nu) = &= \frac{\sum_{v \in V, h_A^t \in H_A^t(h^t)} \mu^*(h^t)(v, h_A^t) \beta^{\mathbf{M}_t^C}(\nu, a_i|v)}{\sum_{v' \in V, h_A^t \in H_A^t(h^t)} \mu^*(h^t)(v', h_A^t) \sum_{l=1}^{N^\nu} \beta^{\mathbf{M}_t^C}(\nu, a'|v')} \\
&= \frac{\sum_{v \in V, h_A^t \in H_A^t(h^t)} \mu^*(h^t)(v, h_A^t) \beta^{\mathbf{M}_t^C}(\nu, a_1|v) \frac{k_i^\nu}{k_1^\nu}}{\sum_{v' \in V, h_A^t \in H_A^t(h^t)} \mu^*(h^t)(v', h_A^t) \sum_{l=1}^{N^\nu} \beta^{\mathbf{M}_t^C}(\nu, a_1|v) \frac{k_l^\nu}{k_1^\nu}} \\
&= \frac{k_i^\nu \left( \sum_{v \in V, h_A^t} \mu^*(h^t)(v, h_A^t) \beta^{\mathbf{M}_t^C}(\nu, a_1|v) \right) / k_1^\nu}{\sum_{l=1}^{N^\nu} k_l^\nu \left( \sum_{v' \in V, h_A^t} \mu^*(h^t)(v', h_A^t) \beta^{\mathbf{M}_t^C}(\nu, a_1|v') \right) / k_1^\nu} \\
&= \frac{k_i^\nu}{\sum_{l=1}^{N^\nu} k_l^\nu}.
\end{aligned}$$

Moreover,

$$\beta^{\mathbf{M}_t'^C}(\nu|v) = \sum_{l=1}^{N^\nu} \beta^{\mathbf{M}_t^C}(\nu, a_l|v) = \sum_{l=1}^{N^\nu} \beta^{\mathbf{M}_t^C}(\nu, a_1|v) \frac{k_l^\nu}{k_1^\nu} = \frac{\beta^{\mathbf{M}_t^C}(\nu, a_1|v)}{k_1^\nu} \sum_{l=1}^{N^\nu} k_l^\nu.$$

Then,

$$\beta^{\mathbf{M}_t'^C}(\nu|v) \sum_{l=1}^{N^\nu} \alpha^{\mathbf{M}_t'^C}(a_l|\nu) \mathbb{E}^{\Gamma^*,\pi^*,r^*}\left[U(a(h^t),a_l,\cdot,v)|h_A^t,\mathbf{M}_t^C,1,\nu,a_l\right] =$$

$$\left(\frac{\beta^{\mathbf{M}_t^C}(\nu,a_1|v)}{k_1^\nu} \sum_{l=1}^{N^\nu} k_l^\nu\right) \times \sum_{l=1}^{N^\nu} \frac{k_l^\nu}{\sum_{j=1}^{N^\nu} k_j^\nu} \mathbb{E}^{\Gamma^*,\pi^*,r^*}\left[U(a(h^t),a_l,\cdot,v)|h_A^t,\mathbf{M}_t^C,1,\nu,a_l\right]$$

$$= \sum_{l=1}^{N^\nu} \beta^{\mathbf{M}_t^C}(\nu,a_1|v) \frac{k_l^\nu}{k_1^\nu} \mathbb{E}^{\Gamma^*,\pi^*,r^*}\left[U(a(h^t),a_l,\cdot,v)|h_A^t,\mathbf{M}_t^C,1,\nu,a_l\right]$$

$$= \sum_{l=1}^{N^\nu} \beta^{\mathbf{M}_t^C}(\nu,a_l|v) \mathbb{E}^{\Gamma^*,\pi^*,r^*}\left[U(a(h^t),a_l,\cdot,v)|h_A^t,\mathbf{M}_t^C,1,\nu,a_l\right].$$

Hence, we have not modified the agent's payoffs.

## III   PROOF OF COROLLARY A.1

In the Appendix, we claim that it follows from Proposition A.1 that for any PBE assessment, $\langle\Gamma^*,(\pi_v^*,r_v^*)_{v\in V},\mu^*\rangle$, there is a payoff-equivalent PBE assessment where the principal's beliefs at each public history $h^t$ satisfy that

$$\mu^*(h^t)(v,h_A^t) = \mu^*(h^t)(v,h_A'^t)$$

whenever $h_A^t, h_A^t{}' \in H_A^t(h^t)$ and $\mu^*(h^t)(v,h_A^t), \mu^*(h^t)(v,h_A^t{}') > 0$. Clearly, Proposition A.1 implies that this holds on the equilibrium path. We now show that the same can be done off the equilibrium path.

Thus, consider a PBE assessment, $\langle\Gamma^*,(\pi_v^*,r_v^*)_{v\in V},\mu^*\rangle$, such that the agent's equilibrium strategy only depends on her payoff-relevant type and the public history. Let $h^t$ be the shortest length public history off the equilibrium path that satisfies that there exists $v \in V$, $h_A^t, h_A^t{}' \in H_A^t(h^t)$ with $\mu^*(h^t)(v,h_A^t), \mu^*(h^t)(v,h_A^t{}') > 0$ and $\mu^*(h^t)(v,h_A^t) \neq \mu^*(h^t)(v,h_A^t{}')$.

Because the agent's strategy does not condition on the payoff-irrelevant part of her private history, the principal's beliefs at histories $h^\tau$ on the path of the equilibrium strategy profile given $h^t$ depend on the agent's payoff-irrelevant private history only through $h_A^t$. That is, if $h_A^\tau, h_A^\tau{}'$ are both successors of $h_A^t$, then $\mu^*(h^\tau)(v,h_A^\tau) = \mu^*(h^\tau)(v,h_A^\tau{}')$.

We now modify the principal's beliefs at history $h^t$ so that they do not depend on the agent's payoff-irrelevant private history. Similar calculations as in the proof of Proposition A.1 then show that the principal's payoff does not change; hence, his strategy remains a best response.

Define $H_A^{t^+}(h^t)(v) = \{h_A^t \in H_A^t(h^t) : \mu^*(h^t)(v, h_A^t) > 0\}$ and let

$$\mu^{**}(h^t)(v, h_A^t) = \begin{cases} \dfrac{\sum_{h_A^t \in H_A^t(h^t)} \mu^*(h^t)(v, h_A^t)}{|H_A^{t^+}(v)|} & \text{if } h_A^t \in H_A^{t^+}(v) \\ 0 & \text{otherwise.} \end{cases}$$

We now modify the principal's beliefs in the continuation histories to reflect the change in the principal's "prior". Fix $\tau \geqslant t+1$. For any history, $h^\tau$, on the path of the equilibrium strategy profile given $h^t$, the principal's beliefs that the agent is of type $v$ and her payoff irrelevant private history is $h_A^{\tau+1} = (h_A^\tau, \mathbf{M}_\tau, 1, m_\tau, s_\tau, a_\tau)$ at history $h^{\tau+1} = (h^\tau, \mathbf{M}_\tau, 1, s_\tau, a_\tau), \mathbf{M}_\tau \in \text{supp } \Gamma^*(h^\tau)$ are given by

$$\mu^{**}(h^{\tau+1})(v, h_A^{\tau+1}) = \frac{\mu^{**}(h^\tau)(v, h_A^\tau) \pi_v^*(h_A^\tau, \mathbf{M}_\tau) r_v^*(h_A^\tau, \mathbf{M}_\tau, 1)(m_\tau) \beta^{\mathbf{M}_\tau}(s_\tau | m_\tau)}{\sum_{\tilde{v}, \widetilde{h_A^\tau}} \sum_{\tilde{m}_\tau} \mu^{**}(h^\tau)(\tilde{v}, \widetilde{h_A^\tau}) \pi_{\tilde{v}}^*(\widetilde{h_A^\tau}, \mathbf{M}_\tau) r_{\tilde{v}}^*(\widetilde{h_A^\tau}, \mathbf{M}_\tau, 1)(\tilde{m}_\tau) \beta^{\mathbf{M}_\tau}(s_\tau | \tilde{m}_\tau)}, \tag{33}$$

and at history $h^{\tau+1} = (h^\tau, \mathbf{M}_\tau, 0, s_\tau, a_\tau), \mathbf{M}_\tau \in \text{supp } \Gamma^*(h^\tau)$ are given by

$$\mu^{**}(h^{\tau+1})(v, h_A^{\tau+1}) = \frac{\mu^{**}(h^\tau)(v, h_A^\tau)(1 - \pi_v^*(h_A^\tau, \mathbf{M}_\tau))}{\sum_{\tilde{v}, \widetilde{h_A^\tau}} \sum_{\tilde{m}_\tau} \mu^{**}(h^\tau)(\tilde{v}, \widetilde{h_A^\tau})(1 - \pi_{\tilde{v}}^*(\widetilde{h_A^\tau}, \mathbf{M}_\tau))}, \tag{34}$$

for $h_A^{\tau+1} \in H_A^{\tau+1}(h^{\tau+1})$.

As in the main text, we now show by induction that for any $\tau \geqslant t$,

$$\sum_{h_A^{\tau+1} \in H_A^{\tau+1}} \mu^{**}(h^{\tau+1})(v, h_A^{\tau+1}) = \sum_{h_A^{\tau+1} \in H_A^{\tau+1}(h^{\tau+1})} \mu^*(h^{\tau+1})(v, h_A^{\tau+1}). \tag{35}$$

As in the main text, we do so for those histories at which the agent participates. It is immediate that this also holds for those histories at which she does not.

For $\tau = t$ and $h^{t+1} = (h^t, \mathbf{M}_t, 1, s_t, a_t)$, the denominator on the right-hand side of

equation ($33$) can be written as:

$$\sum_{\tilde{v},h_A^t} \sum_{m\in M^{\mathbf{M}_t}} \mu^{**}(h^t)(\tilde{v},h_A^t)\pi_{\tilde{v}}^*(h_A^t,\mathbf{M}_t)r_{\tilde{v}}^*(h_A^t,\mathbf{M}_t,1)(m)\beta^{\mathbf{M}_t}(s_t|m)$$

$$=\sum_{\tilde{v}} \sum_{m\in M^{\mathbf{M}_t}} \pi_{\tilde{v}}^*(h_A^t,\mathbf{M}_t)r_{\tilde{v}}^*(h_A^t,\mathbf{M}_t,1)(m)\beta^{\mathbf{M}_t}(s_t|m)\sum_{h_A^t\in H_A^{t^+}(h^t)(v)} \mu^{**}(h^t)(\tilde{v},h_A^t)$$

$$=\sum_{\tilde{v},h_A^t} \sum_{m\in M^{\mathbf{M}_t}} \mu^*(h^t)(\tilde{v},h_A^t)\pi_{\tilde{v}}^*(h_A^t,\mathbf{M}_t)r_{\tilde{v}}^*(h_A^t,\mathbf{M}_t,1)(m)\beta^{\mathbf{M}_t}(s_t|m),$$

where the first equality uses that the agent's strategy does not depend on $h_A^t$ and that only $h_A^t \in H_A^{t^+}(v)$ have positive probability and the second equality uses the definition of $\mu^{**}(h^t)(v,h_A^t)$. It then follows that for $h^{t+1} = (h^t,\mathbf{M}_t,1,s_t,a_t)$ and $h_A^{t+1} = (h_A^t,\mathbf{M}_t,1,m_t,s_t,a_t)$

$$\sum_{h_A^t}\mu^{**}(h^{t+1})(v,h_A^{t+1}) = \frac{\sum_{h_A^t}\mu^*(h^t)(v,h_A^t)\pi_v^*(h_A^t,\mathbf{M}_t)r_v^*(h_A^t,\mathbf{M}_t,1)(m_t)\beta^{\mathbf{M}_t}(s_t|m_t)}{\sum_{\tilde{v},h_A^t}\sum_{m\in M^{\mathbf{M}_t}}\mu^{**}(h^t)(\tilde{v},h_A^t)\pi_{\tilde{v}}^*(h_A^t,\mathbf{M}_t)r_{\tilde{v}}^*(h_A^t,\mathbf{M}_t,1)(m)\beta^{\mathbf{M}_t}(s_t|m)}$$

$$=\sum_{h_A^t}\frac{\mu^*(h^t)(v,h_A^t)\pi_v^*(h_A^t,\mathbf{M}_t)r_v^*(h_A^t,\mathbf{M}_t,1)(m_t)\beta^{\mathbf{M}_t}(s_t|m_t)}{\sum_{\tilde{v},h_A^t}\sum_{m\in M^{\mathbf{M}_t}}\mu^*(h^t)(\tilde{v},h_A^t)\pi_{\tilde{v}}^*(h_A^t,\mathbf{M}_t)r_{\tilde{v}}^*(h_A^t,\mathbf{M}_t,1)(m)\beta^{\mathbf{M}_t}(s_t|m)}$$

$$=\sum_{h_A^t}\mu^*(h^{t+1})(v,h_A^{t+1}).$$

Adding up both sides over $m_t \in M^{\mathbf{M}_t}$ delivers equation ($35$) for $\tau = t$. Similar steps to those in the proof of Proposition A.1 show that indeed equation ($35$) holds for $\tau \geqslant t+1$. As in that proof, this is enough to show that the principal's payoff does not change under the new beliefs. Thus, his strategy remains a best response.

## IV   Continuum of types

### IV.1   Preliminaries

**Primitives:** Let $V,A$ denote compact, metrizable spaces. $V$ denotes the set of agent types, endowed with a full support distribution $\mu_0$. $A$ denotes the set of allocations. Let $M,S$ denote two Polish spaces. All sets are endowed with their Borel $\sigma-$ algebras.

There is a set of *ex-post* allocations, $Y$, also compact and metrizable. The set of feasible

ex-post allocations can depend on the allocation and is captured by a correspondence $\mathcal{Y}$, a measurable subset of $A \times Y$. For each allocation $a \in A$, $\mathcal{Y}_a = \{y \in Y : (a, y) \in \mathcal{Y}\}$.

Given a Polish set $Z$, let $\mathcal{P}(Z)$ denote the set of Borel probability measures on $Z$ endowed with the weak* topology, $\sigma(\mathcal{P}(Z), C_b(Z))$. If $X$ is any other measurable space, a *transition probability* is a measurable mapping $\gamma : X \mapsto \mathcal{P}(Z)$. That is, for any Borel set $Z' \subseteq Z$, $\gamma(Z|x)$ is a measurable function of $x \in X$.

All product sets are endowed with the product topology and the product Borel $\sigma-$algebra.

The principal's von-Neumann Morgenstern utility function is $w : A \times Y \times V \mapsto \mathbb{R}$, whereas the agent's is $u : A \times Y \times V \mapsto \mathbb{R}$. Both functions are measurable.

A mechanism is any $(\beta, \alpha)$ such that $\beta : M \mapsto \mathcal{P}(S)$ and $\alpha : S \mapsto \mathcal{P}(A)$ are transition probabilities. Let $\Gamma$ denote the set of all mechanisms.

**Equilibrium:** A strategy for the principal consists of a choice of mechanism and the specification of a mixed action conditional on every $s \in S$ he may observe. A strategy for the agent maps each of her types to a distribution over messages. To keep matters simple, we avoid discussing the agent's participation decision, but dealing with it is routine.

The focus is on equilibrium outcomes of this game, in which equilibrium means PBE. It consists of a mechanism and three measurable maps: a strategy for $A$, $r : V \mapsto \mathcal{P}(M)$; an ex-post choice for $P$, $\gamma : S \times A \mapsto \mathcal{P}(Y)$ and a belief system $p : S \times A \mapsto \mathcal{P}(V)$ such that

1. $p$ is obtained from $\mu_0$, $r$ and $\beta$ whenever possible $(\star)$

2. $\gamma(s, a)$ is supported on $\arg\max_{y \in \mathcal{Y}_a} \int_V w(a, y, v) dp(v|s)$ for all $(s, a) \in S \times A$

3. $r(v)$ is supported on $\arg\max_{m \in M} \int_{S \times A \times Y} u(a, y, v) d\gamma(y|s, a) d\alpha(a|s) d\beta(s|m)$ for all $v \in V$

$(\star)$ Define for each $m \in M$,

$$(\beta \otimes \alpha)^m(S' \times A') = \int_{S'} \alpha(A'|s) d\beta(s|m) \in \mathcal{P}(S \times A),$$

for each measurable $S' \times A' \subseteq A \times S$.

**Remark IV.1.** Letting $\mu(S' \times A'|m) \equiv (\beta \otimes \alpha)^m(S' \times A')$, note that $((\alpha_s = \alpha(\cdot|s))_{s \in S}, \beta(\cdot|m))$ are a *disintegration* of $\mu(\cdot|m)$ (see, e.g., Balder (2001)). The meaning of this is that for any bounded measurable function $\phi : S \times A \mapsto \mathbb{R}$,

$$\int \phi d\mu(\cdot|m) = \int_S (\int_A \phi(s, a) d\alpha_s(a)) d\beta(s|m).$$

Then, Bayesian updating

$$\int_{V'} \int_M (\beta \otimes \alpha)^m(S' \times A') dr(m|v) d\mu_0(v) = \int_V \int_M \int_{S' \times A'} p(V'|s, a) d(\beta \otimes \alpha)^m(s, a) dr(\cdot|v) d\mu_0(v).$$

### IV.2  Revelation Principle

**Canonical Messages:** The agent's reporting strategy and the mechanism $\beta$, induce a measure on $M \times S$, $(r \otimes \beta)^v$, as follows:

$$(r \otimes \beta)^v(M' \times S') = \int_{M'} \beta(S'|m) dr(m|v),$$

for any measurable $M' \times S' \subseteq M \times S$.[45] Note this defines a new transition $\beta^* : V \times \mathcal{P}(S)$ and that $((\beta_m)_{m \in M}, r(\cdot|v))$ are a disintegration of $(r \otimes \beta)^v$.

When her value is $v$, the agent's payoff is given by

$$\int_M \int_S \int_A \int_Y u(a, y, v) d\gamma(y|s, a) d\alpha(a|s) d\beta(s|m) dr(m|v) = \tag{36}$$

$$= \int_S \int_A \int_Y u(a, y, v) d\gamma(y|s, a) d\alpha(a|s) d\beta^*(s|v), \tag{37}$$

where the equality follows from applying backwards the definition of a disintegration.

Let $u_v^* = \max_{m \in M} \int_S \int_A \int_Y u(a, y, v) d\gamma(y|s, a) d\alpha(a|s) d\beta(s|m)$. By definition of $r(m|v)$,

$$\int_M \int_S \int_A \int_Y u(a, y, v) d\gamma(y|s, a) d\alpha(a|s) d\beta(s|m) dr(m|v) = \int_M u_v^* dr(m|v) = u_v^*. \tag{38}$$

Moreover, for any $m \notin \text{supp } r(\cdot|v)$, $u_v^* \geqslant \int_S \int_A \int_Y u(a, y, v) d\gamma(y|s, a) d\alpha(a|s) d\beta(s|m)$. We

---

[45]Recall we endow all product sets with their product Borel $\sigma$-algebra.

argue that $r^*(v) = \delta(v)$ is an optimal reporting strategy when the mechanism is $(\beta^*, \alpha)$ and the ex-post choice is still $\gamma(\cdot|s,a)$. Equation (36) implies $r^*(v) = \delta(v)$ achieves $u_v^*$. Toward a contradiction, suppose that $v' \in V$, $v' \neq v$ exist such that

$$u_v^* = \int_S \int_A \int_Y u(a,y,v) d\gamma(y|s,a) d\alpha(a|s) d\beta^*(s|v) < \int_S \int_A \int_Y u(a,y,v) d\gamma(y|s,a) d\alpha(a|s) d\beta^*(s|v').$$

$$(39)$$

The right-hand side of the above expression equals

$$\int_S \int_A \int_Y u(a,y,v) d\gamma(y|s,a) d\alpha(a|s) d\beta^*(s|v') =$$
$$\int_M \int_S \int_A \int_Y u(a,y,v) d\gamma(y|s,a) d\alpha(a|s) d\beta(s|m) dr(m|v').$$

Therefore, equations (38) and (39) imply

$$\int_M \int_S \int_A \int_Y u(a,y,v) d\gamma(y|s,a) d\alpha(a|s) d\beta(s|m)(dr(m|v) - dr(m|v')) < 0.$$

This implies that a $\epsilon > 0$ and a set $M' \in \mathcal{B}(M) : r(M'|v') > 0$ exist where

$$\int_S \int_A \int_Y u(a,y,v) d\gamma(y|s,a) d\alpha(a|s) d\beta(s|m') > u_v^* + \epsilon$$

for all $m' \in M'$. This contradicts the optimality of $r(\cdot|v)$.

We now check that we have not changed: (a) $P$'s payoff and (b) Bayesian updating.

To see that the principal's payoff is the same as before, notice that

$$\int_V \left( \int_M \int_S \int_A \int_Y w(a,y,v) d\gamma(y|s,a) d\alpha(a|s) d\beta(s|m) dr(m|v) \right) d\mu_0(v) =$$
$$= \int_V \left( \int_S \int_A \int_Y w(a,y,v) d\gamma(y|s,a) d\alpha(a|s) d\beta^*(s|v) \right) d\mu_0(v),$$

by definition of disintegration.

To see that Bayesian updating has not changed, notice that for any measurable $S' \times A' \subseteq$

$S \times A$,

$$\int_{V'} (\beta^* \otimes A')^v (S' \times A') d\mu_0(v) = \int_{V'} \left( \int_{S'} \alpha(A'|s) d\beta^*(s|v) \right) d\mu_0(v)$$

$$= \int_{V'} \left( \int_M \left( \int_{S'} \alpha(A'|s) d\beta(s|m) \right) dr(m|v) \right) d\mu_0(v)$$

$$= \int_V \int_M \int_{S' \times A'} p(V'|s,a) d(\beta \otimes \alpha)^m(s,a) dr(\cdot|v) d\mu_0(v)$$

$$= \int_V \int_{S' \times A'} p(V'|s,a) d(\beta^* \otimes \alpha)^v(s,a) d\mu_0(v),$$

which yields the desired expression.

**Canonical Outputs:** Given the result in the previous section, we focus hereafter on mechanisms $(\beta^*, \alpha)$, where $\beta^* : V \mapsto \mathcal{P}(S)$ is a transition probability. For now, focus on mechanisms such that $\alpha : S \mapsto \mathcal{P}(A)$ satisfies that $\alpha(\cdot|s) = \delta_a$ for every $s \in S$. Then Bayesian updating reduces to

$$\int_{V'} \beta^*(S'|v) d\mu_0(v) = \int_V \int_{S'} p(V'|s) d\beta^*(s|v) d\mu_0(v),$$

for every measurable $S' \subseteq S$. Using the notation from before, this says that

$$(\mu_0 \otimes \beta^*)(V' \times S') \equiv \int_V \int_{S'} p(V'|s) d(\mu_0 \otimes \beta^*)(v,s).$$

In what follows, we delineate rigorously how $p : S \mapsto \mathcal{P}(V)$ is constructed. Consider $(V \times S, \mathcal{B}_{V \otimes S})$ and endow it with the measure $(\mu_0 \otimes \beta^*)$. Let the measurable mapping $\pi_{\mathbf{S}} : V \times S \mapsto S$ denote the projection onto $S$ and let $\sigma(\pi_{\mathbf{S}}) = \{\pi_{\mathbf{S}}^{-1}(S') : S' \in \mathcal{B}_S\}$ denote the sigma algebra on $V \times S$ induced by $\pi_{\mathbf{S}}$. Note that $\sigma(\pi_{\mathbf{S}}) \subseteq \mathcal{B}_{V \otimes S}$.

Given a function $f \in L^1(V \times S, \mathcal{B}_{V \otimes S}, (\mu_0 \otimes \beta^*))$, that is, a random variable, the conditional expectation of $f$ with respect to $\sigma(\pi_{\mathbf{S}})$ is the function $\mathbb{E}[f|\sigma(\pi_{\mathbf{S}})]$ that satisfies that for every $D \in \sigma(\pi_{\mathbf{S}})$,

$$\int_D f d(\mu_0 \otimes \beta^*) = \int_D \mathbb{E}[f|\sigma(\pi_{\mathbf{S}})] d(\mu_0 \otimes \beta^*).$$

Recall that $\mathbb{E}[f|\sigma(\pi_{\mathbf{S}})]$ exists and is uniquely defined.

Define the random variable $X(v, s) = \mathbb{1}_{V' \times S}(v, s)$ and $\mathbb{P}(V'|S) = \mathbf{E}[X|\sigma(\pi_{\mathbf{S}})]$ as the random variable $Y : V \times S \mapsto \mathbb{R}$ that satisfies that for any $D \in \sigma(\pi_{\mathbf{S}})$,

$$\int_D Y d(\mu_0 \otimes \beta^*) = \int_D X d(\mu_0 \otimes \beta^*).$$

In particular, if $D = D' \times D''$,

$$\int_{D'} \int_{D''} \mathbb{P}(V'|S) d\beta^* d\mu_0 = \int_{D'} \int_{D''} \mathbb{1}_{V' \times S}(v, s) d\beta^*(s|v) d\mu_0(v) = \int_{D' \cap V'} \beta^*(D''|v) d\mu_0(v)$$

$$= \int_V \int_{D''} p(D' \cap V'|s) d(\mu_0 \otimes \beta^*)$$

$$= \int_V \int_{D''} p(V'|s) d(\mu_0 \otimes \beta^*)$$

where the next to last equality follows from the definition of $p$ by Bayesian updating and the last equality follows from noting that if $D = D' \times D'' \in \sigma(\pi_{\mathbf{S}})$, then $D = V \times D'', D'' \in \mathcal{B}_S$. Thus, $p(V'|s)$ is a version of the conditional probability $\mathbf{P}(V'|S)$.

Letting $\mathcal{P}(V)$ denote the space of probability measures on $(V, \mathcal{B}(V))$, we can let $p : V \times S \mapsto \mathcal{P}(V)$ be defined by $p_{(v,s)}(\cdot) = \mathbb{P}(\cdot|\sigma(\pi_{\mathbf{S}}))(v, s)$. It follows that $\int_{V \times S} p_{v,s} d(\mu_0 \times \beta^*) = \mu_0$.

Assume finally that for each $p, a$, a unique maximizer $y^*$ exists. Define $W : \{(p, a, \gamma) : p \in \mathcal{P}(V), a \in A, \gamma(Y') > 0 \Rightarrow Y' \subseteq \mathcal{Y}_a\} \mapsto \mathbb{R}$ to be

$$W(p, a, \gamma) = \int_V \int_Y w(a, y, v) d\gamma(y) dp(v).$$

Then, by the definition of conditional expectation,

$$\int_V \left( \int_S \int_Y w(a(s), y, v) d\gamma(y|s, a) d\beta^*(s|v) \right) d\mu_0(v)$$

$$= \int_V \int_S \mathbb{E}[\int_Y w(a(s), y, v) d\gamma(y|s, a) | \sigma(\pi_{\mathbf{S}})] d\beta^*(s|v) d\mu_0(v)$$

$$= \int_V \int_S W(p_{v,s}, a(s), \gamma) d(\mu_0 \times \beta^*)(s, v)$$

$$= \int_{\mathcal{P}(V)} W(p, a, \gamma) d\tau(p),$$

where the last equality follows from defining for any Borel set $P$ of $\mathcal{P}(V)$, the measure $\tau(P) = (\mu_0 \otimes \beta^*)(\{(v,s) : p_{v,s} \in P\})$.

## IV.3   The envelope representation of payoffs

An advantage of the continuum assumption is that incentive compatibility of the mechanisms implies the agent's payoff from the mechanism can be represented via the envelope theorem. We now derive the corresponding version for our setting. This implies that in the program analyzed in Section 4, downward-looking incentive constraints are always binding, so the relaxed program provides the correct benchmark.

We now make assumptions on the agent's utility function so that we can apply the envelope theorem of Milgrom and Segal (2002). Assume $V$ is a compact subset of the real line. Assume the agent's utility function $u(a, y, \cdot)$ is Lipschitz continuous and that an integrable function $b : V \mapsto \mathbb{R}$ existssuch that $|u_v(a, y, v)| \leqslant b(v)$. We first show this condition implies

$$U(\hat{v}, \cdot) = \int_S \int_A \int_Y u(a, y, v) d\gamma(y|s, a) d\alpha(a|s) d\beta(s|\hat{v})$$

is absolutely continuous and differentiable for all $\hat{v} \in V$.

To see this, note that taking $v, v'$,

$$
\begin{aligned}
|U(\hat{v}, v) - U(\hat{v}, v')| &= |\int_S \int_A \int_Y (u(a, y, v) - u(a, y, v')) d\gamma(y|s, a) d\alpha(a|s) d\beta(s|\hat{v})| \\
&\leqslant \int_S \int_A \int_Y |u(a, y, v) - u(a, y, v')| d\gamma(y|s, a) d\alpha(a|s) d\beta(s|\hat{v}) \\
&\leqslant K|v - v'|,
\end{aligned}
$$

where the last inequality follows from Lipschitz continuity of $u(a, y, \cdot)$. That $U(\hat{v}, \cdot)$ is differentiable follows from the application of Lebesgue's dominated convergence theorem as $v' \to v$ in the above expression.

Then, defining

$$\mathbf{U}(v) = \max_{\hat{v} \in V} U(\hat{v}, v),$$

Theorem 2 in Milgrom and Segal (2002) implies $\mathbf{U}$ is differentiable almost everywhere and

$$\mathbf{U}(v) = \mathbf{U}(\underline{v}) + \int_{\underline{v}}^{v} U_v(t, v)dt,$$

where $U_v(t, v) = \int_A \int_S \int_Y u_v(a, y, t)d\gamma(y|s, a)d\alpha(a|s)d\beta(s|t)$. In particular, in the environment in Section 4 where $A = Q \times \mathbb{R}$, it follows that

$$\begin{aligned}
\mathbb{U}(\underline{v}) + \int_S t_s d\beta(s|v) &= \int_S \int_Q \int_Y u(q, y, v)d\gamma(y|s, a)d\alpha(q|s)d\beta(s|v) \\
&\quad - \int_{\underline{v}}^{v} \int_A \int_S \int_Y u_v(a, y, t)d\gamma(y|s, a)d\alpha(a|s)d\beta(s|t)dt.
\end{aligned}$$