

# A Reputation for Honesty\*

Drew Fudenberg<sup>†</sup>

Ying Gao<sup>‡</sup>

Harry Pei<sup>§</sup>

November 2, 2020

We analyze situations in which players build reputations for honesty rather than for playing particular actions. A patient player facing a sequence of short-run opponents makes an announcement about their intended action after observing an idiosyncratic shock, and before players act. The patient player is either an honest type whose action coincides with their announcement, or an opportunistic type who can freely choose their actions. We show that the patient player can secure a high payoff by building a reputation for being honest when the short-run players face uncertainty about which of the patient player's actions are currently feasible, but may receive a low payoff when there is no such uncertainty.

Many economic actors have reputations for keeping or breaking their promises. As a prominent example, Archibald Cox Jr's default on his promise of paying high bonuses in the early 90s triggered a massive defection of key personnel from First Boston to its archrival Merrill Lynch.<sup>1</sup> Similar logic applies to advertising and marketing, which can set customers' expectations about the types of interactions they are going to have with the firm. If those expectations are not aligned with the actual customer experience, the firm's brand and business will suffer.

Motivated by these observations, we examine the reward for building a reputation for honesty. Compared to reputations for taking specific actions, a reputation for honesty can better adapt an agent's decisions to the current circumstances, which is valuable when the environment changes over time. Moreover, it is unrealistic to make commitments based on future contingencies that are hard to describe in advance, and the simplicity of a commitment to honesty makes it more plausible.

---

\*We thank Mehmet Ekmekci and Navin Kartik for helpful comments, and National Science Foundation grants SES-1947021 and SES-1951056 for financial support.

<sup>†</sup>Department of Economics, Massachusetts Institute of Technology. Email: drewf@mit.edu

<sup>‡</sup>Department of Economics, Massachusetts Institute of Technology. Email: yinggao@gmail.com

<sup>§</sup>Department of Economics, Northwestern University. Email: harrydp@northwestern.edu

<sup>1</sup>See "Taking the Dare" in *The New Yorker*, July 26th, 1993. The departed individuals include leaders of First Boston's prestigious energy group, and more than a dozen managing directors in its fixed-income and mortgage-backed-securities groups, triggered by a lower bonus payment than what they had been promised. Many who departed had been at First Boston their entire careers, including during its difficult times in the 80's.

In our model, a patient player (e.g., a firm) faces a sequence of myopic opponents (e.g., consumers), each of whom plays the game only once. Each period, before players act, the patient player privately observes an idiosyncratic shock, which can affect their payoff (e.g., their production cost) and which of their actions are currently feasible. Then the patient player announces the action they intend to play. The myopic players cannot observe the shocks, but can observe the announcement in the current period as well as whether the patient player has kept their word in the past.

The patient player is either an *honest type*, who strategically chooses their announcements but always keeps their word, or an *opportunistic type*, who strategically chooses both the announcements and the actions. Both types have the same payoff function. This contrasts to Kreps and Wilson (1982), Milgrom and Roberts (1982), and Fudenberg and Levine (1989) in which with positive probability, the patient player is a commitment type who mechanically plays a particular action.

Theorem 1 shows that the patient player receives at least their expected Stackelberg payoff in every equilibrium when the myopic players face a small amount of uncertainty about the actions currently available to the patient player.<sup>2</sup> A complication is that the opportunistic type may announce certain actions with higher probability than the honest type does, so the patient player's announcement may adversely affect their opponent's belief about their type. As a result, both types of the patient player may face a tradeoff between announcing actions that lead to higher credibility and announcing actions that lead to higher commitment payoffs (i.e., payoff conditional on being trusted).

To see why the reputation bound nevertheless obtains, suppose the honest type announces their Stackelberg action whenever it is feasible. When a myopic player does not best reply against the announcement, whether the patient player keeps their word in that period is informative about their type. Because the set of feasible actions is stochastic, the honest type announces each action with strictly positive probability, which implies that observing the current announcement leads to at most a bounded change in the myopic player's belief. Therefore, when the patient player behaves honestly, there can be at most a bounded number of periods in which the myopic players do not best reply to the announcement. As a result, the patient player receives at least their expected Stackelberg payoff.

By contrast, Theorem 2 shows that when the patient player can choose from any of their possible actions in every period, there are equilibria in which they receive a low payoff, which can be as low

---

<sup>2</sup>In Section 4, we show that our reputation result extends when the patient player observes which of their actions are feasible after making their announcement, or when the patient player chooses an action (e.g., their effort), observes their product quality, and makes an announcement about quality before the myopic player chooses their action.

as their minmax value in examples such as the product choice game.

**Related Literature:** Our paper contributes to the study of reputation models where no types are committed to specific actions. Schmidt (1993) characterizes the Markov equilibria of finite-horizon repeated bargaining games in which a firm has private information about its production cost. Pei (2020) characterizes an informed player’s highest Nash equilibrium payoff when facing uninformed opponents. Sugaya and Wolitzky (2020) constructs a cooperative equilibrium in a community enforcement model with a type that communicates strategically but is committed to playing *always defect*. By contrast, we provide a lower bound on the patient player’s payoff for all Nash equilibria.

Our reputation result requires the uninformed players to face uncertainty about the availability of the informed player’s actions, or more generally, believe that the honest type makes every announcement with positive probability. This is related to Celentani, Fudenberg, Levine, and Pesendorfer (1996) and Atakan and Ekmekci (2015), which show that full support monitoring can help reputation building when the uninformed player is long-lived. Their results, unlike ours, require that the informed player cannot perfectly observe the uninformed player’s actions.

Jullien and Park (2020) studies repeated buyer-seller games in which a seller privately observes their product quality, which is a noisy signal of their effort. It shows that cheap talk communication about quality improves the maximum social welfare if and only if the seller’s cost of effort is intermediate.<sup>3</sup> Our paper examines a complementary question, namely, whether a patient player can guarantee high payoffs in *all* equilibria by building reputations for honesty. Successful reputation building in our model requires uncertainty about the actions available to the patient player, but does not depend on the players’ payoff functions. Corollary 2 in Section 4 extends our insights to Jullien and Park (2020)’s setting, which implies that a patient seller receives their optimal commitment payoff in all equilibria when product quality (i.e., the seller’s private signal) is a noisy signal of effort, but receives a payoff lower than that in some equilibria when quality is a perfect signal of effort.

The fact that many people prefer to be honest has been established experimentally by e.g. Gneezy (2005) and Gneezy, Kajackaite, and Sobel (2018). Kartik, Ottaviani, and Squintani (2007) and Kartik (2009) show how costs of lying change the equilibrium outcomes of strategic communication games.

---

<sup>3</sup>Jullien and Park (2014) shows that communication accelerates consumer learning when product quality is determined by the seller’s type, and the high type seller is non-strategic and always tells the truth. Awaya and Krishna (2016) identifies a class of games in which players can achieve perfectly collusive payoffs with communication, but not without it.

Instead of positing that some players have a cost of lying, we follow Chen, Kartik, and Sobel (2008) and Chen (2011) and assume that the patient player is either an honest type who never lies, or an opportunistic type who faces no cost of lying. Our results extend to cases with strictly positive and possibly heterogeneous lying costs.

Our work is related to the literature on pre-play communication. Including an honest type in our model is in line with the experimental finding of Charness and Dufwenberg (2006) that some people keep their word in order to live up to others' expectations. Sobel (2017) allows one of the players to communicate their intended action before playing a two-player complete information game, and provides sufficient conditions under which the sender receives their highest Nash equilibrium payoff.

## 1 Example: Product Choice Game with Stochastic Cost

Consider a game between a firm (row player) with discount factor  $\delta \in (0, 1)$  and a sequence of consumers (column player), each of whom plays the game only once. In every period, the firm privately observes its i.i.d. cost of production  $\theta_t \in \{\theta_g, \theta_b\}$ . Let  $p_g \in (0, 1)$  be the probability that  $\theta_t = \theta_g$ . The players' stage-game payoffs are:

$\theta = \theta_g$	$T$	$N$	$\theta = \theta_b$	$T$	$N$
$H$	1, 2	-1, 0	$H$	-1, 2	-3, 0
$L$	2, -2	0, 0	$L$	2, -2	0, 0

When  $\theta = \theta_g$ , the best pure-strategy commitment for the firm is to action  $H$ , which yields payoff 1.<sup>4</sup> When  $\theta = \theta_b$ , the firm's optimal commitment action is  $L$ , which yields payoff 0. If the firm obtains its optimal pure-strategy commitment payoff in every state, then its expected payoff is  $p_g$ .

**No Announcement Benchmark:** Suppose the firm cannot make announcements about its intended actions, and that with small but positive probability it is a commitment type that mechanically plays  $H$  in every period. Future consumers can observe the firm's effort in previous periods, but not the past realizations of  $\theta_t$ .<sup>5</sup> Then, there are equilibria in which the patient firm's payoff is  $\max\{0, 2p_g - 1\}$ , which is strictly lower than  $p_g$ . For example, when  $p_g \geq 1/2$ , there is an equilibrium where the

<sup>4</sup>A commitment to a mixed strategy is even better in state  $\theta_g$ . We do not consider reputations for playing mixed actions in this paper.

<sup>5</sup>This is a reasonable assumption given that  $\theta$  only affects the firm's cost of supplying high quality.

opportunistic firm chooses  $H$  in every period on the equilibrium path, and each consumer chooses  $T$  unless they observe  $L$  in at least one of the previous periods. Intuitively, when  $\theta = \theta_b$ , the cost of playing  $H$  outweighs the benefit from the consumer's trust, and the firm faces a tradeoff between sustaining its reputation for playing  $H$  and avoiding the excessive cost.

This low-payoff equilibrium motivates our interest in reputations for honesty.

**Reputation for Honesty:** Suppose that the firm can make an announcement  $m_t$  about its intended action  $a_t$  to the current consumer after observing  $\theta_t$ , but before players choosing actions.

The firm is either honest or opportunistic. In contrast to the commitment types in canonical reputation models, the honest type is strategic when making announcements and does not commit to any particular action. Instead, it commits to play the action it announces in every period. The two types of the firm have the same stage-game payoff function and discount factor, and maximize their respective discounted average payoffs. The consumer in period  $t$  observes the firm's announcement in period  $t$ , as well as the value of  $\mathbf{1}\{a_s = m_s\}$  for  $s \in \{0, 1, \dots, t-1\}$ , i.e., whether the firm's announcements matched its actions in the previous periods.

As Theorem 2 shows, the firm's equilibrium payoff can be low when all of its actions are always available. To see how this works in the example, consider the following strategy profile: Both types of the firm announce  $L$  and play  $L$  at every history, and each consumer plays  $N$  regardless of the firm's announcement. The consumers' belief about the firm's type never changes on the equilibrium path. After the firm announces  $H$ , the current consumer believes that the firm is opportunistic and will play  $L$ .<sup>6</sup> This strategy profile and assessment constitute a Perfect Bayesian equilibrium, in which the firm's discounted average payoff is 0 regardless of its type.

This low-payoff equilibrium is driven by the honest-type firm's strategic concerns when making announcements. The consumers believe that the opportunistic type is more likely to announce  $H$ , so the honest type faces a trade-off in state  $\theta_g$  between announcing an action that leads to higher credibility (i.e., action  $L$ ) and an action that leads to a higher commitment payoff (i.e., action  $H$ ). This motivates the honest type to announce  $L$ , making consumers' beliefs self-fulfilling.

In contrast, Theorem 1 shows that when some of the firm's actions are unavailable with small

---

<sup>6</sup>When future consumers only observe whether  $a_t$  coincides with  $m_t$ , but not the exact realizations of  $a_t$  and  $m_t$ , they do not observe deviations in the announcement stage if the firm kept its word. We show that the firm can also receive a low payoff when future consumers can observe both  $a_t$  and  $m_t$ .

but positive probability, both types of the firm receive at least their expected Stackelberg payoff in every equilibrium. For example, suppose in every period the firm can choose between  $H$  and  $L$  with probability  $1 - 2\varepsilon$ , can only choose  $H$  with probability  $\varepsilon$ , and can only choose  $L$  with probability  $\varepsilon$ . Both types of the patient firm receive payoff at least  $(1 - \varepsilon)p_g - 7\varepsilon(1 - p_g)$  when the feasibility of actions is i.i.d. over time and is independent of  $\theta$ . This guaranteed payoff converges to  $p_g$  as  $\varepsilon \rightarrow 0$ .

**Remarks:** Our assumption that the consumers face uncertainty about which of the firm's actions are feasible fits situations in which the firm is a single contractor who occasionally may be sick, and so unable to provide high-quality service. It also fits cases where the firm faces occasional regulatory inspections, and producing low-quality products in those periods can lead to fines and the risk of being shut down. In this situation, the firm will always choose to supply high quality, regardless of their discount factor.<sup>7</sup>

Our reputation result (Theorem 1) extends to cases where the distribution of the feasible actions varies exogenously over time or is correlated with the current  $\theta$ . It also extends when  $\theta_t$  is drawn from a potentially different set  $\Theta_t$  in every period, as long as the patient player's payoff is uniformly bounded for all  $t \in \mathbb{N}$ . This captures situations in which the client's demand varies over time and which of the firm's action benefits the client is known only after the client arrives.

In these situations, the advantage of establishing a reputation for honesty is more pronounced. Due to the complicated nature of future payoff environments, it is impractical for the firm to commit to state-contingent action plans. A reputation for honesty allows the firm to communicate its intended actions after observing the payoff environment, which sidesteps these complications.

## 2 Baseline Model

Time is discrete, indexed by  $t = 0, 1, \dots$ . A long-lived player 1 (e.g., a seller) with discount factor  $\delta$  interacts with an infinite sequence of short-lived player 2s (e.g., consumers), with  $2_t$  denoting the short-lived player in period  $t$ . Player 1's action set is  $A$  and player 2's action set is  $B$ .

Each period consists of an announcement stage and an action stage. In period  $t$ , an i.i.d. random

---

<sup>7</sup>Formally, if the firm chooses  $L$  when it is inspected, it faces a fine  $f > 0$  and a probability  $q \in (0, 1)$  of shutting down. One can show that there exists  $\underline{f} > 0$  and  $\underline{q} \in (0, 1)$  such that when  $f > \underline{f}$  and  $q > \underline{q}$ , it is a dominant strategy for both types of the firm to choose  $a_t = H$  regardless of their discount factors and the equilibrium being played.

variable  $(\theta_t, \omega_t) \in \Theta \times \Omega$  is drawn according to  $p \in \Delta(\Theta \times \Omega)$ , where  $\theta_t \in \Theta$  affects player 1's stage-game payoff (e.g., their cost of supplying high quality) and  $\omega_t \subset A$  is the set of feasible actions, with  $\Omega \equiv 2^A \setminus \{\emptyset\}$ . Player 1 privately observes  $(\theta_t, \omega_t)$  and announces to player 2<sub>t</sub> that they intend to play action  $m_t \in A$ . Players then simultaneously choose their actions  $a_t \in \omega_t$  and  $b_t \in B$ .

Player 1's stage-game payoff is  $u_1(\theta_t, a_t, b_t)$  and player 2<sub>t</sub>'s is  $u_2(a_t, b_t)$ . Each player 2 who arrives after period  $t$  can observe  $y_t \in Y$ , distributed according to  $F(\cdot | a_t, m_t)$ . A leading example is when  $y_t$  is the indicator function  $\mathbf{1}\{a_t = m_t\}$ , that is, future short-run players observe whether the patient player has kept their word in the past.

Player 1 has private information about their type  $\gamma \in \{\gamma_h, \gamma_o\}$ , which is either *honest* ( $\gamma_h$ ) or *opportunistic* ( $\gamma_o$ ). Both types share the same stage-game payoff function. The honest type is restricted (i) to announce an action that is currently available, i.e.,  $m_t \in \omega_t$ , and (ii) to take an action that matches their announcement, i.e.,  $a_t = m_t$ . The opportunistic type can announce any action (including ones that are not feasible that period) and can take any action in  $\omega_t$  regardless of their announcement. Let  $\pi_0 \in (0, 1)$  be the prior probability of the honest type according to player 2s' prior belief.

For every  $t \in \mathbb{N}$ , player 2<sub>t</sub>'s private history is  $h_2^t \equiv \{y_0, y_1, \dots, y_{t-1}, m_t\}$ , with  $h_2^t \in \mathcal{H}_2^t$ . Player 2<sub>t</sub>'s strategy is  $\sigma_2^t : \mathcal{H}_2^t \rightarrow \Delta(B)$ , with  $\sigma_2 \equiv (\sigma_2^t)_{t \in \mathbb{N}}$ . Player 1's private history in the announcement stage of period  $t$  is

$$\widehat{h}_1^t \equiv \{\theta_s, \omega_s, m_s, a_s, b_s, y_s\}_{s=0}^{t-1} \cup \{\gamma, \omega_t, \theta_t\},$$

with  $\widehat{h}_1^t \in \widehat{\mathcal{H}}_1^t$  and  $\widehat{\mathcal{H}}_1 \equiv \bigcup_{t=0}^{\infty} \widehat{\mathcal{H}}_1^t$ . Player 1's private history in the action stage of period  $t$  is

$$\widetilde{h}_1^t \equiv \{\theta_s, \omega_s, m_s, a_s, b_s, y_s\}_{s=0}^{t-1} \cup \{\gamma, \omega_t, \theta_t, m_t\},$$

with  $\widetilde{h}_1^t \in \widetilde{\mathcal{H}}_1^t$  and  $\widetilde{\mathcal{H}}_1 \equiv \bigcup_{t=0}^{\infty} \widetilde{\mathcal{H}}_1^t$ . The opportunistic type's strategy is  $\sigma_o \equiv (\widehat{\sigma}_o, \widetilde{\sigma}_o)$ , with  $\widehat{\sigma}_o : \widehat{\mathcal{H}}_1 \rightarrow \Delta(A)$  their strategy to make announcements and  $\widetilde{\sigma}_o : \widetilde{\mathcal{H}}_1 \rightarrow \Delta(A)$  their strategy to take actions, subject to a feasibility constraint that the support of  $\widetilde{\sigma}_o(\widetilde{h}_1^t)$  is a subset of  $\omega_t$ . The honest type's strategy is  $\sigma_h \equiv (\widehat{\sigma}_h, \widetilde{\sigma}_h)$ , with  $\widehat{\sigma}_h : \widehat{\mathcal{H}}_1 \rightarrow \Delta(A)$  their strategy to make announcements and  $\widetilde{\sigma}_h : \widetilde{\mathcal{H}}_1 \rightarrow \Delta(A)$  their strategy to take actions, subject to first, the support of  $\widehat{\sigma}_h(\widehat{h}_1^t)$  is a subset of  $\omega_t$ , and second, their action matches their announcement  $\widetilde{\sigma}_h(\widetilde{h}_1^t) = m_t$ .

A Nash equilibrium (NE) consists of  $(\sigma_o, \sigma_h, \sigma_2)$ , in which  $\sigma_2^t$  maximizes player 2<sub>t</sub>'s stage-game payoff, and every type of player 1 chooses a strategy that maximizes their discounted average payoff

$\mathbb{E} \left[ \sum_{t=0}^{\infty} (1 - \delta) \delta^t u_1(\theta_t, a_t, b_t) \right]$ . We assume that  $\Theta$ ,  $A$ ,  $B$ , and  $Y$  are finite sets, which together with discounting of per period payoffs implies that a Nash equilibrium exists (Fudenberg and Levine 1983).

### 3 Results

We show that when the short-run players face a small amount of uncertainty about the feasibility of the patient player's actions, and  $y_t$  is informative about whether the patient player has kept their word in period  $t$ , the patient player can secure their expected (pure) Stackelberg payoff in every equilibrium.<sup>8</sup> By contrast, the patient player receives a low payoff in some equilibria when all of their actions are feasible in every period.

Recall that  $\omega_t \subset A$  is the set of feasible actions in period  $t$ . For every  $\varepsilon > 0$ , we say that player 1's action choice is  $\varepsilon$ -flexible if the probability with which  $\omega_t = A$  is at least  $1 - \varepsilon$ .

**Assumption 1.** For every  $a \in A$ ,  $\omega_t = \{a\}$  with strictly positive probability.

Our next assumption requires  $y_t$  to be informative about whether player 1's action and announcement match. A leading example that satisfies this assumption is  $y_t = \mathbf{1}\{a_t = m_t\}$ .

**Assumption 2.** If  $a = m$  and  $a' = m'$ , then (i)  $F(\cdot|a, m) = F(\cdot|a', m')$ , and (ii)  $F(\cdot|a, m)$  does not belong to the convex hull of  $\{F(\cdot|a', m')\}_{a' \neq m'}$ .

Let  $BR_2 : \Delta(A) \rightarrow 2^B \setminus \{\emptyset\}$  be player 2's best reply correspondence. In state  $\theta \in \Theta$ , player 1's Stackelberg payoff is

$$v_1^*(\theta) \equiv \max_{a \in A} \left\{ \min_{b \in BR_2(a)} u_1(\theta, a, b) \right\},$$

and their expected Stackelberg payoff is  $v_1^* \equiv \sum_{\theta \in \Theta} p(\theta) v_1^*(\theta)$ .

**Theorem 1.** Suppose the environment satisfies Assumptions 1 and 2. For every  $\eta > 0$ , there exist  $\underline{\delta} \in (0, 1)$  and  $\varepsilon > 0$  such that when  $\delta > \underline{\delta}$  and player 1's action choice is  $\varepsilon$ -flexible, each type of player 1 receives payoff at least  $v_1^* - \eta$  in every Nash equilibrium.<sup>9</sup>

<sup>8</sup>In what follows, we will simply say *Stackelberg action* and *Stackelberg payoff*, with "pure" left implicit.

<sup>9</sup>When  $y_t = \mathbf{1}\{a_t = m_t\}$ , a patient player 1 can guarantee payoff approximately  $v_1^*$  in every weak rationalizable outcome defined in Watson (1993) in the perturbed game where player 1 is honest with positive probability.



The proof is in Appendix A. For some intuition, consider the example in which  $y_t = \mathbf{1}\{a_t = m_t\}$ . Let  $a^* : \Theta \rightarrow A$  be any mapping such that  $a^*(\theta) \in \arg \max_{a \in A} \left\{ \min_{b \in \text{BR}_2(a)} u_1(\theta, a, b) \right\}$  for every  $\theta$ . Fix any equilibrium, and consider the honest type's payoff when they announce  $a^*(\theta_t)$  in period  $t$  whenever  $a^*(\theta_t) \in \omega_t$ . The second part of Assumption 2 implies that whether player 1's action coincides with their announcement is informative about their type in the "bad" periods where player 2 fails to best reply to the announcement. Assumption 1 requires that for each  $a \in A$ ,  $\omega_t = \{a\}$  with positive probability, which implies that the honest type makes each announcement with positive probability in every period.<sup>10</sup> The first part of Assumption 2 implies that future myopic players cannot identify the patient player's past announcements if the patient player has never lied. Similar to Fudenberg and Levine (1989), there are at most a bounded number of bad periods in which player 2 does not best reply against the announced action. This implies that both types of the patient player obtain at least their expected Stackelberg payoff.

Theorem 1 requires that the patient player can choose any action in  $A$  with probability close to 1. In fact, as long as the environment satisfies Assumptions 1 and 2, one can establish a reputation result using the same argument as the proof of Theorem 1 after adjusting the definition of expected Stackelberg payoff. For every  $(\theta, \omega) \in \Theta \times \Omega$ , let

$$u_1^*(\theta, \omega) \equiv \max_{a \in \omega} \left\{ \min_{b \in \text{BR}_2(a)} u_1(\theta, a, b) \right\}, \quad (3.1)$$

and let

$$u_1^* \equiv \sum_{(\theta, \omega) \in \Theta \times \Omega} p(\theta, \omega) u_1^*(\theta, \omega).$$

One can show that when the environment satisfies Assumptions 1 and 2, a patient player can guarantee payoff  $u_1^*$  in all equilibria when  $\delta$  is close enough to 1.

Next, Assumption 2 rules out situations in which player 2s can perfectly observe player 1's past announcements, or more generally can observe signals that statistically identify the content of player 1's past announcements even when  $a_t = m_t$ . To study such situations, we extend our result when player 2s can observe informative signals  $z_t$  about the past realizations of  $a_t$  and  $m_t$ , as long as each

---

<sup>10</sup>Our reputation result extends when  $\omega_t$  is the set of feasible announcements. When the honest type trembles and makes each announcement with positive probability, our result also extends to settings where player 1 makes their announcement *before* knowing which of their actions are feasible. See Corollary 3 for details.

of them can only observe the realizations of  $z_t$  in a bounded number of previous periods.

Formally, let  $z_t \in Z$ , where  $z_t$  is distributed according to  $G(\cdot|m_t, a_t)$ . We make no restrictions on  $G$  except that its support  $Z$  is a finite set. Suppose for every  $t \in \mathbb{N}$ , player 2<sub>*t*</sub> can observe player 1's announcement  $m_t$ , the history  $\{y_0, \dots, y_{t-1}\}$ , as well as a (possibly stochastic) subset of  $\{z_0, \dots, z_{t-1}\}$  that has at most  $K \in \mathbb{N}$  elements. Whether player 1 can observe  $y$  and  $z$  are irrelevant for our result.

Our assumption on the asymmetry between player 2s' observations of  $y_t$  and  $z_t$  is motivated by retail markets in developing economies. Due to the lack-of record-keeping institutions, detailed information about sellers' actions and announcements (e.g., the quality of their services, various attributes of their products, the content of their advertisements, and so on, which correspond to  $z_t$ ) is likely to get lost over time. By contrast, coarse information about sellers' records, such as whether they have kept their word (which corresponds to  $y_t$ ), is likely to be more persistent due to social learning and word-of-mouth communication. Corollary 1 extends Theorem 1, with proof in Online Appendix A.

**Corollary 1.** *Suppose the environment satisfies Assumptions 1 and 2, and there exists  $K \in \mathbb{N}$  such that each player 2 observes the past realizations of  $z$  in at most  $K$  periods. For every  $\eta > 0$ , there exist  $\underline{\delta} \in (0, 1)$  and  $\varepsilon > 0$  such that when  $\delta > \underline{\delta}$  and player 1's action choice is  $\varepsilon$ -flexible, each type of player 1 has payoff at least  $v_1^* - \eta$  in every Nash equilibrium.*

Now we show why uncertainty about which of player 1's actions are feasible is necessary for Theorem 1 to hold in general. We show this for situations in which  $\omega_t = A$  with probability 1 and  $y_t = \mathbf{1}\{a_t = m_t\}$ .<sup>11</sup>

We start by introducing two auxiliary one-shot games that have the same payoff functions as the original stage game. The first auxiliary game does not have a communication stage: Player 1 observes  $\theta$ , and then players act simultaneously without any communication. Let  $v_1^{min}$  be player 1's lowest Nash equilibrium payoff in this game. The second auxiliary game has an action recommendation stage: Player 1 observes  $\theta$ , makes a recommendation  $\hat{b} \in B$  to player 2 before players take their actions. Let  $\hat{v}_1$  be player 1's lowest pure-strategy equilibrium payoff in this game. If there is no pure-strategy equilibrium in this game, let  $\hat{v}_1 = +\infty$ .

Let  $\mathcal{B}$  be the set of mappings  $\beta : A \rightarrow \Delta(B)$  such that  $\beta(a)$  is a best reply to  $a$  for every  $a \in A$ .

---

<sup>11</sup>We do not know how to construct low-payoff equilibria without the assumption that  $\omega_t = A$  with probability 1 and  $y_t = \mathbf{1}\{a_t = m_t\}$ , and we do not know of any counterexamples.

Abusing notation, let  $p$  be the distribution of  $\theta$ . Let

$$v'_1 \equiv \min_{A' \subset A, \beta \in \mathcal{B}} \sum_{\theta \in \Theta} p(\theta) \max_{a \in A'} u_1(\theta, a, \beta(a)) \quad (3.2)$$

subject to

$$\sum_{\theta \in \Theta} p(\theta) \max_{a \in A'} u_1(\theta, a, \beta(a)) \geq \min\{v_1^{\min}, \widehat{v}_1\}. \quad (3.3)$$

Theorem 2 shows that when all of player 1's actions are feasible in every period, there are equilibria in which both types of player 1 have payoff no more than  $v'_1$ ,

**Theorem 2.** *If  $\omega_t = \{A\}$  with probability 1 and  $y_t = \mathbf{1}\{a_t = m_t\}$ , then there exists  $\underline{\delta} \in (0, 1)$  such that for every  $\delta > \underline{\delta}$ , there exists an equilibrium in which both types of player 1's payoff is  $v'_1$ .*

The proof of this result and a subsequent lemma are in Appendix B.

In order to understand the connections between the conclusion of Theorem 2 and that of Theorem 1, we compare  $v'_1$  with player 1's expected Stackelberg payoff  $v_1^*$  and their minmax payoff. To start with, one can verify that  $v_1^* \geq v'_1$  when players' payoffs are generic and the auxiliary game without communication admits a pure-strategy equilibrium.<sup>12</sup> Next, we introduce a class of games under which  $v'_1$  is strictly less than  $v_1^*$ , and under an additional supermodularity condition,  $v'_1$  equals player 1's minmax payoff.

**Supermodularity Condition.** *There exists a complete order on  $A$  such that for every  $\theta \in \Theta$ ,  $u_1(\theta, a, b)$  is strictly decreasing in  $a$ , and there exists  $\theta \in \Theta$  such that player 1's Stackelberg action in state  $\theta$  is not the lowest element in  $A$ .*

**Lemma 1.** *If every  $\theta \in \Theta$  occurs with positive probability and the stage-game payoffs satisfy supermodularity, then*

1.  $v'_1 < v_1^*$ .
2. *In addition, if there also exists a complete order on  $B$  such that  $u_1$  is strictly increasing in  $b$  and  $u_2$  has strictly increasing differences in  $a$  and  $b$ , then  $v'_1$  is player 1's minmax payoff.*

---

<sup>12</sup>The generic requirement is that player 1 has a strict best reply to every  $b \in B$  for every  $\theta \in \Theta$ , and player 2 has a strict best reply to every  $a \in A$ . The existence of a pure-strategy equilibrium in the auxiliary game without communication rules out zero-sum games such as matching pennies, where the patient player cannot benefit from committing to pure actions.

Condition 1 and the additional assumption in Lemma 1 fit applications such as product choice games, where a firm finds it costly to exert high effort, can strictly benefit from consumers' trust, and can benefit from committing to high effort in states where its production cost is low enough. The consumers have stronger incentives to trust the firm when the latter exerts higher effort. Our conditions also apply to games of entry deterrence (Kreps and Wilson 1982, Milgrom and Roberts 1982), capital taxation (Phelan 2006), monetary policy (Barro 1986), and trust games more generally (Liu and Skrzypacz 2014).

## 4 Extensions

We note here that the conclusion of Theorem 1 extends to two alternative scenarios.

**Announcing Product Quality:** Suppose that players move sequentially in the stage game. In period  $t \in \mathbb{N}$ , player 1 (e.g., a firm) chooses their effort  $a_t \in A$ , privately observes the quality of its product  $x_t \in X$  which is distributed according to  $g(\cdot|a_t) \in \Delta(A)$ , and makes an announcement  $m_t \in X$  about quality. Player 2<sub>*t*</sub> (e.g., a consumer) observes  $m_t$  as well as whether  $x_s$  coincides with  $m_s$  for all  $s \leq t - 1$  before choosing  $b_t \in B$ . We assume that  $A$ ,  $B$ , and  $X$  are finite sets.

Player 1 is either an honest type who strategically chooses actions  $a_t \in A$  but announces  $x_t$  truthfully, or an opportunistic type who strategically chooses both the actions and the announcements. Both types have stage-game payoff  $u_1(a_t, b_t)$  and discount factor  $\delta \in (0, 1)$ . Player 2<sub>*t*</sub>'s payoff is  $u_2(x_t, b_t)$ , i.e., their payoff depends only on product quality and their purchasing decision. This fits the model of Jullien and Park (2020) except that there is a positive probability of the honest type, and the ex-post quality is not directly observed by subsequent consumers. For every  $x \in X$ , let  $\text{BR}_2(x) \subset B$  be the set of pure best replies against  $x$ . Player 1's *optimal commitment payoff* is

$$v^{**} \equiv \max_{a \in A} \left\{ \sum_{x \in X} g(x|a) \min_{b \in \text{BR}_2(x)} u_1(a, b) \right\}. \quad (4.1)$$

**Corollary 2.** *If  $g(\cdot|a)$  has full support for every  $a \in A$ , then for every  $\varepsilon > 0$ , there exists  $\underline{\delta} \in (0, 1)$  such that when  $\delta > \underline{\delta}$ , every type of player 1's payoff in every Nash equilibrium is at least  $v^{**} - \varepsilon$ .*

The proof is in Online Appendix B, which uses similar ideas as the proof of Theorem 1.

Next we show that reputations for honesty cannot guarantee player 1 their optimal commitment payoff when product quality is a perfect signal of player 1's effort. Suppose  $X = A$  and  $g(a|a) = 1$  for every  $a \in A$ , and players' stage-game payoffs are given by the following matrix:

–	$T$	$N$
$H$	1,2	–1,0
$L$	2,–2	0,0

Player 1's optimal commitment payoff is 1, which can be obtained by committing to play  $H$ .

We construct a Perfect Bayesian equilibrium in which player 1's payoff is 0. On the equilibrium path, both types of player 1 play  $L$  and announce  $L$  in every period, and player 2s play  $N$  at every on-path history. After observing announcement  $H$ , player 2s believe that player 1 is the opportunistic type and has played  $L$  with probability  $1/2$  in the current period, and best reply by playing  $N$ . This equilibrium survives both when player 2s can only observe whether  $m_t$  matches with  $x_t$  in all previous periods, and when player 2s can observe the values of  $x_t$  and  $m_t$  in all previous periods.

**Making Announcements Before Knowing the Set of Feasible Actions:** In some applications, the patient player makes announcements before knowing which of their actions are feasible, and an honest individual may break their word when the action they announced turns out to be infeasible. Theorem 1 extends to this setting if (1) the honest type trembles and makes each announcement with positive probability, and (2) the probability with which all actions are feasible is close to 1.

In period  $t \in \mathbb{N}$ , player 1 observes  $\theta_t \in \Theta_t$  and makes an announcement about their intended action  $m_t \in A_t$ . Player 2 <sub>$t$</sub>  observes  $m_t$ , player 1 observes the realization of  $\omega_t \in \Omega \equiv 2^A \setminus \{\emptyset\}$ , and then both players choose  $(a_t, b_t) \in \omega_t \times B$  simultaneously. Future player 2s observe  $y_t \equiv \mathbf{1}\{a_t = m_t\}$ . We assume  $\{\omega_t, \theta_t\}_{t \in \mathbb{N}}$  are i.i.d. over time, with  $p \in \Delta(\Omega \times \Theta)$  their joint distribution.

Player 1 is either an opportunistic type who can take any action regardless of their announcement, or an honest type who chooses  $a_t = m_t$  as long as  $m_t \in \omega_t$ . Both types of player 1 tremble when making announcements, i.e., there exists  $\eta > 0$  such that the probability with which each type makes each announcement is at least  $\eta$  at every information set.

**Corollary 3.** *For every  $\varepsilon > 0$ , there exist  $\underline{\delta} \in (0, 1)$ ,  $\bar{\eta} > 0$ , and  $c \in (0, 1)$ , such that when  $\delta > \underline{\delta}$ ,  $\eta \in (0, \bar{\eta})$  and the probability that  $\omega_t = A$  is at least  $1 - \eta c$ , then each type of player 1 receives payoff at least  $v_1^* - \varepsilon$  in every Nash equilibrium.*

The proof is in Online Appendix C. Unlike our baseline model and reputation models with noisy monitoring such as Fudenberg and Levine (1992), when the honest type uses the strategy of announcing their Stackelberg action and keeping their word whenever it is feasible, their reputation may deteriorate in expectation.

Our proof starts by showing that when  $\omega_t = A$  with probability close to 1, the probability that the honest type keeps their word in equilibrium is close to 1, and for reputation to deteriorate when the honest type keeps their word, the opportunistic type must also keep their word with probability close to 1. It implies that in those periods, player 2 has a strict incentive to best reply to player 1's announcement, and moreover, the amount of reputation deterioration is small. By contrast, in "bad" periods where player 2 has a strict incentive not to best reply against player 1's announcement, the probability that the opportunistic type breaks their promise is large and keeping one's word leads to a significant improvement in one's reputation. Although the number of bad periods can be unbounded, their fraction goes to zero as the probability of  $\omega_t = A$  goes to one.

## A Proof of Theorem 1

Fix any Nash equilibrium  $(\sigma_o, \sigma_h, \sigma_2)$  and consider any history  $h^t$  that occurs with strictly positive probability under  $(\sigma_o, \sigma_h, \sigma_2)$ .

For  $i \in \{h, o\}$ , let  $P^{\sigma_i, \sigma_2}$  be the probability measure over  $Y^\infty$  induced by  $(\sigma_i, \sigma_2)$ . Denote player 2's belief over player 1's private history as a function of  $h_2^t$  by  $\beta(\hat{h}_1^t | h_2^t)$ , and let  $\sigma_o(h_2^t)$  be the expected distribution of opportunistic player 1's joint announcement-action pairs implied by  $\beta(h_1^t | h_2^t)$ , with  $\hat{\sigma}_o(h_2^t)$  and  $\tilde{\sigma}_o(h_2^t)$  the marginal distributions of announcements and actions, respectively. Let  $\pi_t$  be the probability of the honest type according to player 2's belief in period  $t$  after observing  $\{y_0, \dots, y_{t-1}\}$ . According to Bayes rule,

$$\pi_t = \frac{P^{\sigma_h, \sigma_2}(y_0, \dots, y_{t-1})\pi_0}{P^{\sigma_h, \sigma_2}(y_0, \dots, y_{t-1})\pi_0 + P^{\sigma_o, \sigma_2}(y_0, \dots, y_{t-1})(1 - \pi_0)}. \quad (\text{A.1})$$

Let

$$\alpha_t(m_t) \equiv \pi_t \hat{\sigma}_h(h_2^t)(m_t) + (1 - \pi_t) \tilde{\sigma}_o(h_2^t)(m_t),$$

which is the probability of announcement  $m_t$  conditional on  $h_2^t$ . Let  $\xi_t(m_t)$  be the probability that

$a_t = m_t$  conditional on  $m_t$ ,

$$\xi_t(m_t) \equiv \frac{\pi_t \tilde{\sigma}_h(h_2^t)(m_t)}{\pi_t \hat{\sigma}_h(h_2^t)(m_t) + (1 - \pi_t) \hat{\sigma}_o(h_2^t)(m_t)} + \frac{(1 - \pi_t) \tilde{\sigma}_o(h_2^t)(m_t)}{\pi_t \hat{\sigma}_h(h_2^t)(m_t) + (1 - \pi_t) \hat{\sigma}_o(h_2^t)(m_t)}.$$

Let  $\xi_t$  be the unconditional probability that player 1's action matches their announcement:

$$\xi_t \equiv \sum_{a \in A} \alpha_t(a) \xi_t(a). \quad (\text{A.2})$$

Let  $\underline{\rho} \equiv \min_{a \in A} \Pr(\omega_t = \{a\})$ , which by Assumption 1 is strictly positive. This implies that  $\alpha_t(m) > \underline{\rho}$  for every announcement  $m \in A$ . Let  $\bar{\lambda} \in (0, 1)$  be the smallest real number such that for every  $\theta \in \Theta$ , player 2 strictly prefers one of the actions in  $\text{BR}_2(a^*(\theta))$  to all actions outside of  $\text{BR}_2(a^*(\theta))$  when they believe that player 1 plays  $a^*(\theta)$  with probability strictly more than  $\bar{\lambda}$ .

Consider the honest type's payoff when they use strategy  $\sigma_h^* \equiv (\hat{\sigma}_h^*, \tilde{\sigma}_h^*)$ , where  $\tilde{\sigma}_h^*(m) = m$  for every  $m \in A$ , and  $\hat{\sigma}_h^*(\theta_t, \omega_t) = a^*(\theta_t)$  when  $a^*(\theta_t) \in \omega_t$  and is uniform over the actions in  $\omega_t$  when  $a^*(\theta_t) \notin \omega_t$ . For any history  $h^t$ , suppose there exists  $m_t \in A$  such that  $\xi_t(m_t) \leq \bar{\lambda}$ , then  $\xi_t \leq \xi^* \equiv 1 - (1 - \bar{\lambda})\underline{\rho}$ . Let  $d(\cdot || \cdot)$  denote the KL-divergence, and let  $F^* \equiv F(\cdot | a, a)$ . Let

$$D^* \equiv \min_{a \neq m} d\left(F^* \left\| \xi^* F^* + (1 - \xi^*) F(\cdot | a, m)\right.\right). \quad (\text{A.3})$$

Part 2 of Assumption 2 implies that  $D^* > 0$  and is independent of player 1's discount factor  $\delta$ .

Part 1 of Assumption 2 implies that  $P^{\sigma_h, \sigma_2} = P^{\sigma_h^*, \sigma_2}$ . Let  $F(y | h_2^t) = \sum_{(a, m) \in A^2} F(y | a, m) \sigma_o(h_2^t)(a, m)$  so that  $F(\cdot | h_2^t)$  denotes the distribution over  $y_t$  induced by  $\sigma_o(h_2^t)$ .

Similar to Gossner (2011), the chain rule for relative entropy implies:

$$\begin{aligned} -\log \pi_0 &\geq d\left(P^{\sigma_h, \sigma_2} \left\| \pi_0 P^{\sigma_h, \sigma_2} + (1 - \pi_0) P^{\sigma_o, \sigma_2}\right.\right) \\ &= d\left(P^{\sigma_h^*, \sigma_2} \left\| \pi_0 P^{\sigma_h, \sigma_2} + (1 - \pi_0) P^{\sigma_o, \sigma_2}\right.\right) = \mathbb{E}^{(\sigma_h, \sigma_2)} \left[ \sum_{t=0}^{\infty} d\left(F^* \left\| \pi_t F^* + (1 - \pi_t) F(\cdot | h_2^t)\right.\right) \right]. \end{aligned} \quad (\text{A.4})$$

Therefore,  $d(F^* || F(h^t)) \geq D^*$  if  $h^t$  is such that  $\xi_t(a_t) \leq \bar{\lambda}$  for some  $a_t \in A$ , so the expected number of such periods is at most

$$\bar{T}(\pi_0) \equiv \left\lceil \frac{-\log \pi_0}{D^*} \right\rceil \quad (\text{A.5})$$

Hence the honest type's payoff from  $\sigma_h^*$  is at least

$$\delta^{T(\pi_0)} \left\{ \left(1 - \frac{\varepsilon}{\min_{\theta \in \Theta} p(\theta)}\right) v_1^* + \frac{\varepsilon}{\min_{\theta \in \Theta} p(\theta)} \underline{v}_1 \right\} + (1 - \delta^{T(\pi_0)}) \underline{v}_1, \quad (\text{A.6})$$

in which  $\underline{v}_1$  is player 1's lowest stage-game payoff. Expression (A.6) converges to  $v_1^*$  when  $\delta \rightarrow 1$  and  $\varepsilon \rightarrow 0$ . Since the opportunistic type's payoff is weakly greater than the honest type's payoff, their equilibrium payoff is also weakly more than (A.6).

## B Proofs of Theorem 2 and Lemma 1

**Proof of Theorem 2:** Suppose  $(A', \beta)$  solves (3.2) subject to (3.3), and consider the following strategy profile: At every on-path history with  $\theta_t = \theta$ , both types of player 1 announce the same  $a \in \arg \max_{a \in A'} u_1(\theta, a, \beta(a))$  and match their actions with their announcements. Player 2 chooses  $\beta(a)$  following announcement  $a$ , and chooses  $\beta(a')$  if player 1's announcement does not belong to  $A'$ , where  $a'$  is an arbitrary element of  $A'$ . At every  $h^t$  where  $y_s \neq 1$  for some  $s < t$ , player 2 believes that player 1 is opportunistic. If  $v_1^{\min} \leq \widehat{v}_1$ , then players coordinate on the worst stage-game Nash equilibrium for player 1. If  $v_1^{\min} > \widehat{v}_1$ , then players coordinate on the worst pure-strategy Nash equilibrium in the second auxiliary game. Player 2's incentive constraints are trivially satisfied, and player 1's incentive constraint is implied by (3.3).

**Proof of Lemma 1:** Let  $\underline{a} \equiv \min A$  and let  $\underline{b}$  be player 2's best reply against  $\underline{a}$ . According to (3.2),  $v_1' \leq \sum_{\theta \in \Theta} p(\theta) u_1(\theta, \underline{a}, \underline{b})$ . According to (3.1),  $v^* \geq \sum_{\theta \in \Theta} p(\theta) u_1(\theta, \underline{a}, \underline{b})$ , and the inequality is strict from the second part of Condition 1 and the fact that each  $\theta$  has strictly positive probability.

If  $u_1$  is strictly decreasing in  $b$  and  $u_2$  has strictly increasing differences, then  $u_1(\theta, \underline{a}, \underline{b})$  is player 1's minmax payoff in state  $\theta$ . Since  $v_1' \leq \sum_{\theta \in \Theta} p(\theta) u_1(\theta, \underline{a}, \underline{b})$ ,  $v_1'$  is player 1's minmax value.



## References

- Alp Atakan and Mehmet Ekmekci. Reputation in the long-run with imperfect monitoring. *Journal of Economic Theory*, 157:553–605, 2015.
- Yu Awaya and Vijay Krishna. On communication and collusion. *American Economic Review*, 106(2): 285–315, 2016.
- Robert Barro. Reputation in a model of monetary policy with incomplete information. *Journal of Monetary Economics*, 17(1):3–20, 1986.
- Marco Celentani, Drew Fudenberg, David Levine, and Wolfgang Pesendorfer. Maintaining a reputation against a long-lived opponent. *Econometrica*, 64(3):691–704, 1996.
- Gary Charness and Martin Dufwenberg. Promises and partnership. *Econometrica*, 74(6):1579–1601, 2006.
- Ying Chen. Perturbed communication games with honest senders and naive receivers. *Journal of Economic Theory*, 146(2):401–424, 2011.
- Ying Chen, Navin Kartik, and Joel Sobel. Selecting cheap-talk equilibria. *Econometrica*, 76(1): 117–136, 2008.
- Drew Fudenberg and David Levine. Subgame-perfect equilibria of finite and infinite horizon games. *Journal of Economic Theory*, 31(2):251–268, 1983.
- Drew Fudenberg and David Levine. Reputation and equilibrium selection in games with a patient player. *Econometrica*, 57(4):759–778, 1989.
- Drew Fudenberg and David Levine. Maintaining a reputation when strategies are imperfectly observed. volume 59(3), pages 561–579. 1992.
- Uri Gneezy. Deception: The role of consequences. *American Economic Review*, 95(1):384–394, 2005.
- Uri Gneezy, Agne Kajackaite, and Joel Sobel. Lying aversion and the size of the lie. *American Economic Review*, 108(2):419–453, 2018.
- Olivier Gossner. Simple bounds on the value of a reputation. *Econometrica*, 79(5):1627–1641, 2011.
- Bruno Jullien and In-Uck Park. New, like new, or very good? reputation and credibility. *Review of Economic Studies*, 81(4):1543–1574, 2014.
- Bruno Jullien and In-Uck Park. Communication, feedbacks and repeated moral hazard with short-lived buyers. *Working Paper*, 2020.
- Navin Kartik. Strategic communication with lying costs. *Review of Economic Studies*, 76(4):1359–1395, 2009.

- Navin Kartik, Marco Ottaviani, and Francesco Squintani. Credulity, lies, and costly talk. *Journal of Economic Theory*, 134(1):93–116, 2007.
- David Kreps and Robert Wilson. Reputation and imperfect information. *Journal of Economic Theory*, 27(2):253–279, 1982.
- Qingmin Liu and Andrzej Skrzypacz. Limited records and reputation bubbles. *Journal of Economic Theory*, 151:2–29, 2014.
- Paul Milgrom and John Roberts. Predation, reputation, and entry deterrence. *Journal of Economic Theory*, 27(2):280–312, 1982.
- Harry Pei. Trust and betrayals: Reputational payoffs and behaviors without commitment. *Theoretical Economics*, forthcoming, 2020.
- Christopher Phelan. Public trust and government betrayal. *Journal of Economic Theory*, 130(1): 27–43, 2006.
- Klaus Schmidt. Commitment through incomplete information in a simple repeated bargaining game. *Journal of Economic Theory*, 60:114–139, 1993.
- Joel Sobel. A note on pre-play communication. *Games and Economic Behavior*, 102:477–486, 2017.
- Takuo Sugaya and Alexander Wolitzky. Communication and community enforcement. *Working Paper*, 2020.
- Joel Watson. A “reputation” refinement without equilibrium. *Econometrica*, 61(1):199–205, 1993.