

## B Long-run Agents

### B.1 A Result with long-run Agents

#### B.1.1 Model, Result, and Discussion

Consider a repeated game with a single principal and  $N$  agents with a shared discount factor  $\delta \in [0, 1)$ . In each period, the following stage game is played:

1. Exactly one agent is publicly selected to be active. For each agent  $i \in \{1, \dots, N\}$ , let  $x_{i,t} \in \{0, 1\}$  be the indicator function for agent  $i$  being selected. Let  $\Pr\{x_{i,t} = 1\} = \rho_i$ , where  $\sum_i \rho_i = 1$ .
2. The active agent chooses  $e_t \in \mathbb{R}_+$  and  $\mu_t : \mathbb{R} \rightarrow M$ , which are observed only by the principal and the active agent.
3. The principal and the active agent exchange transfers, with resulting net transfer to the active agent  $s_t \in \mathbb{R}$ . These transfers are observed only by the principal and the active agent.
4. The message  $m_t = \mu_t(s_t)$  is realized and publicly observed.

The principal's and agent  $i$ 's payoffs in each period  $t$  are  $\pi_t = e_t - s_t$  and  $u_{i,t} = x_{i,t}(s_t - c(e_t))$ , respectively, with corresponding expected discounted payoffs  $\Pi_t = \sum_{t'=t}^{\infty} \delta^{t'-t}(1 - \delta)(e_t - s_t)$  and  $U_{i,t} = \sum_{t'=t}^{\infty} \delta^{t'-t}(1 - \delta)x_{i,t}(s_t - c(e_{i,t}))$ . Our solution concept is plain Perfect Bayesian Equilibrium with one additional restriction: at any history  $h^t$  such that agent  $i$  has observed a deviation, we require that  $\mathbb{E}[U_{i,t}|h^t] \geq 0$ . This restriction rules out pathological off-path behavior that might arise from the fact that an agent's beliefs about the history are essentially arbitrary once he observes a deviation.<sup>8</sup> We also restrict attention to equilibria in pure strategies to simplify agents' beliefs on the equilibrium path.

---

<sup>8</sup>This condition is trivially satisfied in any equilibrium that is recursive. It is needed here because this game has private monitoring, which means that equilibria are not necessarily recursive.

**Proposition 10** *Let  $e_i^*$  be the maximum effort attainable in any pure-strategy Perfect Bayesian equilibrium. Letting  $s_i^* \equiv \min\{e_i^*, e^{FB}\} - c(\min\{e_i^*, e^{FB}\})$ ,  $e_1^*, e_2^*, \dots, e_N^*$  must satisfy the system of inequalities*

$$(1 - \delta)c(e_i^*) \leq 2\delta\rho_i s_i^* + \frac{2\rho_i\delta}{1 - (1 - \rho_i)\delta} \sum_{j \neq i} \rho_j s_j^*. \quad (18)$$

It is instructive to compare the right-hand side of (18) to the condition  $c(e^*) \leq 3(H - L)$  from Proposition 7. To translate between settings, note that the total surplus created by the principal's future interactions with agent  $i$  equals  $\delta\rho_i s_i^*$ , which corresponds to  $2(H - L)$  in Proposition 7. In Proposition 7, the principal earns an additional  $H - L$  if she refuses to pay an agent who has shirked. In the game with long-run agents, the principal can be given the *entire* continuation surplus from her relationship with agent  $i$ , which accounts for the second  $\delta\rho_i s_i^*$  in the right-hand side of (18).

The second term on the right-hand side of (18) represents a new force for cooperation that is not present in Proposition 7. Since each agent  $i$  chooses a new communication protocol whenever he is active, he essentially commits to his messages *only until he next interacts with the principal again*. An agent might therefore use his future messages to reveal that he has extorted the principal in equilibrium. However, he cannot do so until the next time that he is active, so this term shrinks to zero as  $\rho_i \rightarrow 0$ .

An immediate corollary of Proposition 10 is that, as the probability that an agent interacts with the principal  $\rho_i$  approaches zero, that agent's maximum equilibrium effort does too. This implication is similar to our main takeaway from Proposition 7: the strength of each agent's bilateral relationship limits the severity of the coordinated punishments available to him. This result relies on the fact that agents can send messages only when they are active. We can interpret this assumption as the natural extension of our commitment assumption to a setting with long-run agents; indeed, a result identical to Proposition 10 would hold if agents could communicate in every period but whenever an agent is active, he commits to a

sequence of messages in each period until he is again active.

### B.1.2 Proof of Proposition 10

For each agent  $j \in \{1, \dots, 2\}$ , let  $e_j^*$  be the maximum effort that can be attained in any period of any equilibrium. Consider a history  $h^t$  right after agent  $i$  is chosen to be the active agent in period  $t$ . Define four different expectations of  $\Pi_{t+1}$  that follow four different outcomes:

1.  $\bar{\Pi}^*$  if no player deviates, with corresponding message  $\bar{m}$ ;
2.  $\underline{\Pi}^*$  if the principal deviates but the active agent does not, with corresponding message  $\underline{m}$ ;
3.  $\bar{\Pi}^{HU}$  if the active agent deviates and  $m_t = \bar{m}$ ;
4.  $\underline{\Pi}^{HU}$  if the active agent deviates and  $m_t = \underline{m}$ .

We identify necessary conditions for effort  $e$  to be attained in equilibrium.

First, the principal must be willing to pay  $s^*$  if the active agent does not deviate, which requires

$$s^* \leq \frac{\delta}{1-\delta} (\bar{\Pi}^* - \underline{\Pi}^*). \quad (19)$$

Second, the active agent  $i$  must be willing to choose effort  $e$  and the equilibrium communication protocol  $\mu$ . Agent  $i$  can always deviate by choosing  $e_t = 0$  and

$$\mu_t = \begin{cases} \underline{m} & s_t < \hat{s} \\ \bar{m} & \text{otherwise} \end{cases}$$

for some  $\hat{s} \geq 0$ . Following this deviation, the principal's unique best response is to pay  $\hat{s}$  so long as

$$-\hat{s} + \frac{\delta}{1-\delta} \bar{\Pi}^{HU} > \frac{\delta}{1-\delta} \underline{\Pi}^{HU},$$

since the principal can earn no less than  $\bar{\Pi}^{HU}$  in the continuation game if  $m_t = \bar{m}$  and no more than  $\underline{\Pi}^{HU}$  if  $m_t = \underline{m}$ . Therefore, agent  $i$  has no profitable deviation of this form only if

$$s^* - c(e) + \frac{\delta}{1-\delta} \bar{U}_i^* \geq \max \left\{ 0, \frac{\delta}{1-\delta} \left( \bar{\Pi}^{HU} - \underline{\Pi}^{HU} \right) \right\}, \quad (20)$$

where  $\bar{U}_i^*$  is the agent's expectations about her continuation payoff at the history that yields principal payoff  $\bar{\Pi}_i^*$ .

Combining (19) and (20) yields the following necessary condition for effort  $e$  to be part of equilibrium:

$$c(e) \leq \frac{\delta}{1-\delta} \left( \bar{U}_i^* + \bar{\Pi}^* - \underline{\Pi}^* \right) - \max \left\{ 0, \frac{\delta}{1-\delta} \left( \bar{\Pi}^{HU} - \underline{\Pi}^{HU} \right) \right\} \quad (21)$$

Our next goal is to connect (19) and (20) by studying the relationship between  $\bar{U}_i^* + \bar{\Pi}^* - \underline{\Pi}^*$  and  $\bar{\Pi}^{HU} - \underline{\Pi}^{HU}$ . We do so by bounding  $\bar{U}_i^* + \bar{\Pi}^* - \bar{\Pi}^{HU}$  from above and  $\underline{\Pi}^* - \underline{\Pi}^{HU}$  from below.

Fix two histories  $h^{t+1}$  and  $\hat{h}^{t+1}$  at the start of period  $t+1$  such that agent  $i$  can distinguish  $h^{t+1}$  from  $\hat{h}^{t+1}$  but no other agents can. For  $t' \geq t+1$ , we will use the notation  $h^{t'}$  and  $\hat{h}^{t'}$  to represent successor histories to  $h^{t+1}$  and  $\hat{h}^{t+1}$ , respectively. At history  $\hat{h}^{t+1}$ , the principal can always play the following strategy:

1. At any history  $\hat{h}^{t'}$  that the active agent believes is consistent with  $h^{t+1}$ , play as in the corresponding successor history to  $h^{t+1}$ ;
2. At any other history, choose  $s_t = 0$ .

Under this strategy, each agent  $j \neq i$  learns that the history is inconsistent with  $h^{t+1}$  only when agent  $i$  sends a message that is inconsistent with play following  $h^{t+1}$ . In a pure-strategy

equilibrium, all agents learn this fact at the same time. For each  $t' \geq t + 1$ , denote

$$\hat{\mathcal{B}}^{t'} = \left\{ \hat{h}^{t'} \mid \text{Agents } j \neq i \text{ learn that the history is inconsistent with } h^{t+1} \text{ in period } t' - 1, \right. \\ \left. \text{but not before} \right\}.$$

Where  $\hat{\mathcal{B}}^\infty$  denotes the event that agents  $j \neq i$  never learn that the history is inconsistent with  $h^{t+1}$ . Note that these events collectively partition the set of histories following  $\hat{h}^{t+1}$ . We can define an analogous collection of sets for the event that agents  $j \neq i$  learn that the history is inconsistent with  $\hat{h}^{t+1}$ . We denote this analogous collection  $\mathcal{B}^{t'}$ .

For each agent  $j \in \{1, \dots, N\}$ , define  $\pi_{j,t} = x_{j,t}(e_t - s_t)$  and  $\pi_{-j,t} = \sum_{k \neq j} x_{k,t}(e_t - s_t)$  as the principal's payoff from agent  $j$  and from all other agents, respectively. Define  $\Pi_{j,t} = \sum_{t'=t}^\infty \delta^{t'-t}(1 - \delta)\pi_{j,t'}$  and  $\Pi_{-j,t} = \sum_{t'=t}^\infty \delta^{t'-t}(1 - \delta)\pi_{-j,t'}$ . Because  $\{\hat{\mathcal{B}}^{\tilde{t}}\}_{\tilde{t}=t+1}^{t'}$  partitions the histories of length  $t'$  following  $\hat{h}^{t+1}$ ,

$$\mathbb{E} \left[ \Pi_{t+1} \mid \hat{h}^{t+1} \right] = \sum_{t'=t+1}^\infty (1 - \delta)\delta^{t'-t-1} \left( \mathbb{E} \left[ \pi_{i,t'} \mid \hat{h}^{t+1} \right] + \sum_{\tilde{t}=t+1}^{t'} \mathbb{E} \left[ \pi_{-i,t'} \mid \hat{h}^{t+1}, \hat{\mathcal{B}}^{\tilde{t}} \right] \Pr \left\{ \hat{\mathcal{B}}^{\tilde{t}} \right\} \right). \quad (22)$$

The right-hand side of (22) is absolutely convergent, so we can rearrange the order of summation to yield

$$\mathbb{E} \left[ \Pi_{t+1} \mid \hat{h}^{t+1} \right] = \sum_{\tilde{t}=t+1}^\infty \left( \sum_{t'=t+1}^{\tilde{t}-1} (1 - \delta)\delta^{t'-t-1} \mathbb{E} \left[ \pi_{i,t'} \mid \hat{h}^{t+1} \right] + \right. \\ \left. \sum_{t'=t+1}^{\tilde{t}-1} (1 - \delta)\delta^{t'-t-1} \mathbb{E} \left[ \pi_{-i,t'} \mid \hat{\mathcal{B}}^{\tilde{t}} \right] + \delta^{\tilde{t}-t-1} \mathbb{E} \left[ \Pi_{-i,\tilde{t}} \mid \hat{\mathcal{B}}^{\tilde{t}} \right] \right) \Pr \left\{ \hat{\mathcal{B}}^{\tilde{t}} \right\}. \quad (23)$$

Under the principal's strategy specified above, the principal and agents  $j \neq i$  act identically until those agents learn of a deviation. Therefore, for any  $t' < \tilde{t}$ ,

$$\mathbb{E} \left[ \pi_{-i,t'} \mid \mathcal{B}^{\tilde{t}} \right] \Pr \left\{ \mathcal{B}^{\tilde{t}} \right\} = \mathbb{E} \left[ \pi_{-i,t'} \mid \hat{\mathcal{B}}^{\tilde{t}} \right] \Pr \left\{ \hat{\mathcal{B}}^{\tilde{t}} \right\}.$$

Moreover, for any  $\tilde{t}$ ,  $\Pr \left\{ \mathcal{B}^{\tilde{t}} \right\} = \Pr \left\{ \hat{\mathcal{B}}^{\tilde{t}} \right\}$ , since any message that distinguish  $h^{t+1}$  from  $\hat{h}^{t+1}$  must also distinguish  $\hat{h}^{t+1}$  from  $h^{t+1}$ .

Now,  $\mathbb{E} \left[ \Pi_{t+1} | \hat{h}^{t+1} \right]$  is bounded below by the principal's payoff from the strategy specified above. Therefore, we can use (23) to bound the difference

$$\begin{aligned} & \mathbb{E} \left[ \Pi_{t+1} | h^{t+1} \right] - \mathbb{E} \left[ \Pi_{t+1} | \hat{h}^{t+1} \right] \leq \\ & \sum_{t'=t+1}^{\infty} \delta^{t'-t-1} (1-\delta) \left( \mathbb{E} \left[ \pi_{i,t'} | h^{t+1} \right] - \mathbb{E} \left[ \pi_{i,t'} | \hat{h}^{t+1} \right] \right) + \\ & \sum_{\tilde{t}=t+1}^{\infty} \delta^{\tilde{t}-t-1} \left( \mathbb{E} \left[ \Pi_{-i,\tilde{t}} | \mathcal{B}^{\tilde{t}} \right] - \mathbb{E} \left[ \Pi_{-i,\tilde{t}} | \hat{\mathcal{B}}^{\tilde{t}} \right] \right) \Pr \left\{ \mathcal{B}^{\tilde{t}} \right\} \end{aligned} \quad (24)$$

Under the specified strategy,  $\mathbb{E} \left[ \Pi_{-i,\tilde{t}} | \hat{\mathcal{B}}^{\tilde{t}} \right] \geq 0$  because the principal pays no transfer to an agent  $j$  who knows that the history is inconsistent with  $h^{t+1}$ , with  $\mathbb{E} \left[ \pi_{i,t'} | \hat{h}^{t+1} \right] \geq 0$  for a similar reason. A necessary condition for (24) is therefore

$$\begin{aligned} & \mathbb{E} \left[ \Pi_{t+1} | h^{t+1} \right] - \mathbb{E} \left[ \Pi_{t+1} | \hat{h}^{t+1} \right] \leq \\ & \sum_{t'=t+1}^{\infty} \delta^{t'-t-1} \left( (1-\delta) \mathbb{E} \left[ \pi_{i,t'} | h^{t+1} \right] + \mathbb{E} \left[ \Pi_{-i,t'} | \mathcal{B}^{t'} \right] \Pr \left\{ \mathcal{B}^{t'} \right\} \right) \end{aligned} \quad (25)$$

Suppose that  $h^{t+1}$  is the on-path history such that  $\mathbb{E} \left[ \Pi_{t+1} | h^{t+1} \right] = \bar{\Pi}^*$ . In a pure-strategy equilibrium, agents correctly infer the true history on the equilibrium path, which means that they must earn nonnegative utility. Consequently, the principal earns no more than total continuation surplus, so (25) requires

$$\mathbb{E} \left[ \Pi_{t+1} | h^{t+1} \right] - \mathbb{E} \left[ \Pi_{t+1} | \hat{h}^{t+1} \right] \leq \sum_{t'=t+1}^{\infty} \delta^{t'-t-1} \left( (1-\delta) \rho_i s_i^* + \sum_{j \neq i} \rho_j s_j^* \Pr \left\{ \mathcal{B}^{t'} \right\} \right). \quad (26)$$

Note that an identical bound holds for the expression  $\bar{U}_i^* + \bar{\Pi}^* - \underline{\Pi}^*$  because agents  $j \neq i$  earn nonnegative continuation utilities on the equilibrium path. If  $h^{t+1}$  is instead the history such that  $\mathbb{E} \left[ \Pi_{t+1} | h^{t+1} \right] = \underline{\Pi}^{HU}$ , then agents have observed  $\underline{m}$  and so know that play is off-path. Our equilibrium restriction requires their utilities to be nonnegative at such a history, so (26) again holds.

Since  $s_j^* \geq 0$ , the right-hand side of (26) is maximized by having the event  $\mathcal{B}^{\tilde{t}}$  happen as

early as possible. The earliest it can occur is the next time that agent  $i$  is the active agent, since agent  $i$  can send a message only when he is active. Agent  $i$  is active for the first time since period  $t$  in period  $t'$  with probability  $(1 - \rho_i)^{t'-t-1}\rho_i$ , so (26) requires

$$\begin{aligned} \mathbb{E} [\Pi_{t+1}|h^{t+1}] - \mathbb{E} [\Pi_{t+1}|\hat{h}^{t+1}] &\leq \sum_{t'=t+1} \delta^{t'-t-1} \left( (1 - \delta)\rho_i s_i^* + (1 - \rho_i)^{t'-t-1}\rho_i \sum_{j=1}^N \rho_j s_j^* \right) \\ &= \rho_i s_i^* + \frac{\rho_i}{1 - (1 - \rho_i)\delta} \sum_{j \neq i} \rho_j s_j^*. \end{aligned} \tag{27}$$

As argued above, an identical bound holds for  $\mathbb{E} [U_{i,t+1} + \Pi_{t+1}|h^{t+1}] - \mathbb{E} [\Pi_{t+1}|\hat{h}^{t+1}]$ .

From (27), we conclude that

$$\bar{U}_i^* + \bar{\Pi}^* - \underline{\Pi}^* \leq \bar{\Pi}^{HU} - \underline{\Pi}^{HU} + 2\rho_i s_i^* + \frac{2\rho_i}{1 - (1 - \rho_i)\delta} \sum_{j \neq i} \rho_j s_j^*.$$

A necessary condition for (21) to hold is therefore

$$c(e) \leq \left( \begin{array}{c} \frac{\delta}{1-\delta} \left( \bar{\Pi}^{HU} - \underline{\Pi}^{HU} + 2\rho_i s_i^* + \frac{2\rho_i}{1-(1-\rho_i)\delta} \sum_{j \neq i} \rho_j s_j^* \right) - \\ \max \left\{ 0, \frac{\delta}{1-\delta} \left( \bar{\Pi}^{HU} - \underline{\Pi}^{HU} \right) \right\} \end{array} \right).$$

The right-hand side of this condition is maximized by  $\bar{\Pi}^{HU} - \underline{\Pi}^{HU} = 0$ , in which case

$$(1 - \delta)c(e) \leq 2\rho_i s_i^* + \frac{2\rho_i}{1 - (1 - \rho_i)\delta} \sum_{j \neq i} \rho_j s_j^*,$$

as desired. ■

## C Communication by the Principal

### C.1 The Principal Can Send Messages

Let  $M_p$  be the set of messages for the principal, and  $m_p$  a typical message. In each period  $t \geq 0$ , the principal chooses a message  $m_{p,t}$  in each period  $t \geq 0$ , and this message is publicly

observed. We consider two different stage games: the principal might either choose  $m_{p,t} \in M_p$  before or after agent  $t$  chooses  $m_t$ . If the principal chooses  $m_{p,t}$  before  $m_t$  is realized, we assume that  $\mu_t$  is a function of  $s_t$  only (and so doesn't depend on  $m_{p,t}$ ).

**The principal talks after agent  $t$ .** Consider some period  $t$ . We let  $\pi(m, m_p)$  be the principal's continuation payoff if  $(m, m_p)$  realizes. Given agent  $t$ 's message  $m$ , the principal always chooses  $m_p$  to maximize  $\pi(m, m_p)$ . We let  $\pi(m) := \max_{m_p} \pi(m, m_p)$ , so  $\pi(m)$  is the principal's continuation payoff after agent  $t$ 's message  $m$ . We let  $\bar{\Pi}$  and  $\underline{\Pi}$  be the highest and lowest continuation payoffs that agent  $t$ 's message can induce. Then, incentive constraints are identical to the the extortion game (i.e., Proposition 2). The principal's message does not mitigate extortion at all, so our impossibility result still holds.

**Proposition 11** *Suppose that in each period  $t$  the principal sends  $m_p \in M_p$  after agent  $t$  sends  $m$ . The principal-optimal equilibrium is outcome-equivalent to that in Proposition 2.*

**The principal talks before agent  $t$ .** Consider some period  $t$ . Define  $\pi(m_p, m)$  as the principal's continuation payoff if  $m_t = m$  and  $m_{p,t} = m_p$ . Once the principal chooses  $s_t$ , she knows  $m_t = \mu_t(s_t)$ . The principal therefore chooses  $m_{p,t}$  to maximize her continuation payoff given agent  $t$ 's message.<sup>9</sup> The same argument as in the previous case applies, so every equilibrium involves zero effort in each period.

## C.2 The Principal Can Commit to a Communication Protocol

In this appendix, we modify the extortion game by allowing the principal to choose a communication protocol at the same time as each agent. We first show that Proposition 1 holds in this game, which means that allowing the principal to commit to messages as a function

---

<sup>9</sup>This intuition would not change if agents could commit to a mixture over  $M$ , in which case the principal would choose  $m_{p,t}$  to maximize her continuation payoff given the mixture. The key is that agent  $t$  can use her message to implement the same punishment regardless of whether he works or shirks.



of transfers eliminates extortion. We then give two reasons why this result should be treated with skepticism.

Formally, suppose that in each  $t \geq 0$ , the principal chooses a communication protocol  $\nu_t : \mathbb{R} \rightarrow M$  at the same time that agent  $t$  chooses  $e_t$  and  $\mu_t$ . At the end of  $t$ , message  $m_t^P = \nu_t(s_t)$  is realized and publicly observed (along with agent  $t$ 's message  $m_t$ ). We can adapt the proof of Proposition 1 to show that the principal can earn no more than  $e^* - c(e^*)$  in this game, where  $e^*$  is defined as in Proposition 1. It suffices to construct an equilibrium in which she earns that payoff.

Consider the following strategy profile. Play starts in the cooperation phase. In this phase,

$$\nu_t(s_t) = \mu_t(s_t) = \begin{cases} C & s_t \geq c(e^*) \\ D & \text{otherwise} \end{cases}$$

and  $e_t = e^*$ . If neither player deviates, then  $s_t = c(e^*)$ ; if only agent  $t$  deviates, then  $s_t = 0$ ; if the principal or both players deviate, then the principal best-responds given the communication protocols. The game stays in the cooperative phase if  $m_t = m_t^P = C$ . Otherwise, it switches to the punishment phase with probability  $\gamma \in [0, 1]$ . In the punishment phase, agents exert no effort and the principal pays no transfers.

Choosing  $\gamma$  to solve

$$c(e^*) = \frac{\delta}{1 - \delta} \gamma (e^* - c(e^*)) \quad (28)$$

implies that the principal is willing to pay  $s_t = c(e^*)$  on the equilibrium path. If agent  $t$  deviates, then the principal's continuation payoff cannot exceed  $e^* - c(e^*)$  if she pays  $s_t = c(e^*)$  and equals  $(1 - \gamma)(e^* - c(e^*))$  if she pays any other amount. Condition (28) implies that she is willing to pay  $s_t = 0$  in that case. Agent  $t$  therefore has no profitable deviation from  $e_t$  or  $\mu_t$ . The principal has no profitable deviation from  $\nu_t$ , since given  $\mu_t$ , she earns no more than  $e^* - c(e^*)$  for paying  $s_t = c(e^*)$  and no more than  $(1 - \gamma)(e^* - c(e^*))$  for paying any other amount. This strategy profile is therefore an equilibrium. It is principal-

optimal because it maximizes total equilibrium surplus and gives all of that surplus to the principal.

This argument shows that allowing the principal to commit to a communication protocol eliminates extortion. Essentially, the principal's and each agent's communication protocols can be used to "cross-check" one another. If the principal is punished whenever messages disagree, then agents cannot extort any *smaller* amount than the amount that the principal pays a hard-working agent on-path. As in the proof of Proposition 5, the principal can then be made indifferent between paying  $s_t = c(e^*)$  and  $s_t = 0$ , so that she is willing to pay a hard-working agent but not one that shirks.

While allowing the principal to commit to a communication protocol can in principle restore cooperation, this result should be treated with skepticism for two reasons. First, while agents are indifferent across messages, the principal is not. Indeed, appendix C.1 shows that she has a strict incentive to send the message that maximizes her continuation payoff. Commitment therefore forces the principal to send messages that she strictly prefers not to send, which stands in contrast to the agents, for whom commitment simply breaks indifference across messages. Consequently, we cannot treat the principal's communication protocol as an equilibrium refinement; no analogue to Proposition 8 exists for the game with principal commitment.

Second, as appendix C.1 illustrates, this result requires the principal to choose  $\nu_t$  (*weakly*) *before* agent  $t$  chooses  $\mu_t$  and  $e_t$ . If agent  $t$  chooses  $\mu_t$  first, then he can shirk and extort the principal, in which case her unique best-response is to pay that agent and then send a message that guarantees a high continuation payoff. If the principal chooses  $\nu_t$  before agent  $t$  chooses  $\mu_t$ , in contrast, then we can slightly modify the equilibrium construction above to show that a version of Proposition 1 holds. The conclusion that principal commitment eliminates extortion therefore depends on a particular assumption about *when* each player makes threats.

## D Variants of the Extortion Game

### D.1 Up-Front Transfers

The **extortion game with up-front transfers** is identical to the extortion game except that at the start of each period  $t$ , the principal and agent  $t$  exchange nonnegative transfers. Denote the resulting net wage to agent  $t$  by  $w_t \in \mathbb{R}$ , so that the principal's and agent  $t$ 's payoffs are  $e_t - s_t - w_t$  and  $s_t + w_t - c(e_t)$ , respectively. Note that agent  $t$  chooses a communication mechanism  $\mu_t$  only *after* these transfers are paid.

**Proposition 12** *Any equilibrium of the extortion game with up-front transfers entails  $e_t = s_t = 0$ .*

#### Proof of Proposition 12

Borrowing notation from the proof of Proposition 2,  $s_t \leq \frac{\delta}{1-\delta}(\bar{\Pi} - \underline{\Pi})$  on the equilibrium path, and any  $s_t < \frac{\delta}{1-\delta}(\bar{\Pi} - \underline{\Pi})$  can be made a unique best response after agent  $t$  deviates. Therefore,  $c(e_t) = 0$  in any  $t$  of any equilibrium. ■

### D.2 *Ex Ante* Extortion

If the principal and each agent can exchange up-front transfers, as they do in appendix D.1, it is natural to consider equilibria if agents can commit to their communication protocols as a function of those transfers. In particular, an agent might demand an up-front transfer in exchange for refraining from later extortionary threats. Of course, once the principal pays an up-front transfer to an agent, that agent has every incentive to renege on her earlier promise not to engage in extortion. In this section, we prove that even if agents can commit to not engage in future extortion in exchange for up-front payments, the unique equilibrium outcome still entails zero effort in each period.

To make this point, consider the follow game. The **game with *ex ante* extortion** is identical to the extortion game with up-front transfers, except that at the start of each

period  $t$ , agent  $t$  chooses a **communication meta-protocol**  $\mu_t^0 : \mathbb{R}^2 \rightarrow \mathcal{M}$ , where

$$\mathcal{M} \equiv \{\mu : \mathbb{R} \rightarrow M\}$$

is the set of communication protocols. This meta-protocol is observed by the principal but not by other agents. The principal and agent  $t$  then exchange up-front transfers, with net transfer  $w_t$ , and agent  $t$  chooses  $e_t \geq 0$ . Agent  $t$ 's communication protocol equals the one specified by the meta-protocol, given  $(w_t, e_t)$ :

$$\mu_t^0(w_t, e_t)(\cdot).$$

The rest of the period proceeds with this communication protocol.

While this alternative game might seem cumbersome, it is designed to capture a very simple intuition. In the extortion game, extortion involves shirking and so is inefficient. In principle, an agent could more efficiently extort the principal by demanding an up-front transfer in exchange for refraining from further extortion. Clearly, it might be difficult for an agent to commit to not make future extortionary threats. Even if he can overcome this commitment problem, however, *ex ante* extortion does not facilitate cooperation. In particular, each agent can use *ex ante* extortion to demand the entire proceeds from his effort. But then the principal earns nothing in any period, which means that she is unwilling to compensate any agent for his effort. The resulting unique equilibrium outcome entails no effort.

**Proposition 13** *Every equilibrium of the game with ex ante extortion entails  $e_t = 0$  in each  $t \geq 0$ .*

### Proof of Proposition 13

Define  $\Pi^* \geq 0$  as the principal's maximum equilibrium payoff, and consider an equilibrium that attains  $\Pi^*$ . If  $e_0 = 0$  with probability 1 in this equilibrium, then

$$\Pi^* \leq \delta \Pi^*$$

and so  $\Pi^* = 0$ .

Suppose that  $e_0 > 0$  with positive probability in this equilibrium. Define  $\bar{\Pi}$  and  $\underline{\Pi}$  as the largest and smallest equilibrium continuation payoffs induced by  $m_0$  in this equilibrium, with corresponding messages  $C$  and  $D$ , respectively. Let  $e^*$  equal the effort that maximizes total period-0 surplus among all on-path efforts in period 0. Then  $e^* > 0$ , and moreover, it must be that

$$\frac{\delta}{1-\delta}(\bar{\Pi} - \underline{\Pi}) \geq c(e^*) > 0.$$

Fix some  $\hat{w} \geq 0$ . For any  $\epsilon, \xi > 0$ , consider the following choice of  $\mu_0^0$  by agent 0:

$$\mu_0^0(w_0, e_0)(\cdot) = \begin{cases} \mu^C(\cdot) & \text{if } e_t = e^* - \epsilon, w_t = \hat{w} \\ \mu^C(\cdot) & \text{if } e_t = 0, w_t \neq \hat{w} \\ \mu^D(\cdot) & \text{otherwise} \end{cases},$$

where

$$\mu^C(s_0) = \begin{cases} C & s_0 = \frac{\delta}{1-\delta}(\bar{\Pi} - \underline{\Pi}) - \xi \\ D & \text{otherwise} \end{cases}$$

and

$$\mu^D(s_0) = D.$$

Give this choice, suppose that  $w_0 \neq \hat{w}$ . If  $e_0 = 0$ , then the principal's unique best

response (to  $\mu^C$ ) is

$$s_0 = \frac{\delta}{1-\delta} (\bar{\Pi} - \underline{\Pi}) - \xi.$$

If  $e_0 > 0$ , then the principal's unique best response (to  $\mu^D$ ) is  $s_0 = 0$ . Therefore, agent  $t$ 's uniquely optimal effort is  $e_0 = 0$ , in which case the principal's payoff is at most

$$(1-\delta)\xi + \delta\underline{\Pi}.$$

Suppose instead that  $w_0 = \hat{w}$ . If  $e_0 = e^* - \epsilon$ , then the principal's unique best response (to  $\mu^C$ ) is again

$$s_0 = \frac{\delta}{1-\delta} (\bar{\Pi} - \underline{\Pi}) - \xi.$$

If  $e_0 \neq e^* - \epsilon$ , then the principal's unique best response (to  $\mu^D$ ) is  $s_0 = 0$ . For any  $\epsilon > 0$ , there exists a sufficiently small  $\xi > 0$  such that

$$-c(e^* - \epsilon) + \frac{\delta}{1-\delta} (\bar{\Pi} - \underline{\Pi}) - \xi > 0.$$

Therefore,  $e_0 = e^* - \epsilon$  is agent 0's uniquely optimal effort, in which case the principal's payoff is

$$(1-\delta)(\xi + e^* - \epsilon - \hat{w}) + \delta\underline{\Pi}.$$

We have uniquely pinned down the principal's payoff as a function of  $\hat{w}$ . The principal's finds it strictly optimal to pay  $w_0 = \hat{w}$  so long as

$$(1-\delta)(\xi + e^* - \epsilon - \hat{w}) + \delta\underline{\Pi} > (1-\delta)\xi + \delta\underline{\Pi}$$

or

$$\hat{w} < e^* - \epsilon.$$

Agent 0 can therefore use this strategy to guarantee a payoff arbitrarily close to

$$e^* - \epsilon - c(e^* - \epsilon) + \frac{\delta}{1 - \delta}(\bar{\Pi} - \underline{\Pi}) - \xi. \quad (29)$$

In equilibrium, agent 0's utility and the principal's payoff cannot exceed

$$(1 - \delta)(e^* - c(e^*)) + \delta\bar{\Pi}.$$

Subtracting (29) from this surplus and taking  $\epsilon, \xi \rightarrow 0$ , we conclude that the principal's payoff cannot exceed  $\delta\underline{\Pi}$ . But then

$$\Pi^* \leq \delta\underline{\Pi} \leq \delta\Pi^*,$$

so again  $\Pi^* = 0$ .

We have established that the principal's *maximum* equilibrium payoff equals 0, which is also her min-max payoff. Therefore,  $\bar{\Pi} = \underline{\Pi} = 0$ , which implies that  $e_t = 0$  in each  $t \geq 0$  of any equilibrium. ■

While every equilibrium entails zero effort in the game with *ex ante* extortion, the intuition for this result differs from that of Proposition 2. Here, each agent can use the meta-communication protocol to extract the entire surplus created by his interaction with the principal. The principal therefore has no reason to actually pay an agent, since her continuation payoff equals zero regardless of her actions today. Proposition 13 is a particularly extreme consequence of the negative intertemporal externality from Proposition 4. In this case, each agent's rent-seeking behavior is so severe that it totally undermines cooperation, resulting in zero effort in equilibrium.

This model assumes that agent  $t$ 's meta-protocol conditions on both the up-front transfer and his effort. We make this assumption to draw a close connection to the extortion model, since agents in that model implicitly condition their communication protocols on their efforts.

Note that assuming  $\mu_t^0$  can condition on  $e_t$  does not resolve the commitment problem at the heart of the model, since agent  $t$  must still find it sequentially optimal to exert effort given his beliefs about  $s_t$ . We could instead assume that  $\mu_t^0$  can condition on  $w_t$  but *not* on  $e_t$ , in which case we can construct equilibria with strictly positive effort. These equilibria take advantage of the fact that once an agent deviates in  $\mu_t^0$ , there might exist a continuation equilibrium in which he expects  $s_t = 0$  and so exerts no effort. Analogous to Proposition 7, the possibility of inefficient continuation play *within* period  $t$  limits agent  $t$ 's ability to engage in *ex ante* extortion. The principal can therefore earn a strictly positive equilibrium payoff if meta-protocols are not contingent on effort.