

# Making Corruption Harder: Asymmetric Information, Collusion, and Crime\*

Juan Ortner

Sylvain Chassang<sup>†</sup>

Boston University

New York University

February 28, 2017

## Abstract

We model the investigation of criminal activity as a principal-agent-monitor problem in which the agent can corrupt the monitor and side-contract to destroy evidence. Building on insights from Laffont and Martimort (1997) we study whether the principal can benefit from endogenously creating asymmetric information between the agent and the monitor. We show that the principal can benefit from randomizing the incentives given to the monitor (and letting those serve as the monitor’s private information), but that the optimality of random incentives depends on pre-existing patterns of private information. We address the issue by providing a data-driven framework for policy evaluation that requires only unverified report data. A potential local policy change is an improvement if, everything else equal, it is associated with greater reports of crime.

KEYWORDS: monitoring, collusion, corruption, asymmetric information, random incentives, prior-free policy evaluation.

---

\*An early version of this paper was circulated under the title “Making Collusion Hard: Asymmetric Information as a Counter-Corruption Measure.” We are especially thankful to Gerard Padró i Miquel for many helpful and inspiring conversations. We are indebted to Yeon-Koo Che, Hugo Hopenhayn, Bart Lipman, Dilip Mookherjee, Stephen Morris, Andy Newman, Debraj Ray as well as seminar participants at Arizona State University, Caltech, the Canadian Economic Theory Conference 2015, Collegio Carlo Alberto, Columbia, the 13th Columbia/Duke/MIT/Northwestern IO Theory Conference, the Einaudi Institute, Iowa State, LSE, McGill, Minnesota, Rice, SITE, the ThReD Conference 2016, UCL, UT Austin, University of Toronto, Warwick, and Wash. U. for feedback. We have also benefited from the thoughtful feedback of a Co-Editor and anonymous referees. Chassang gratefully acknowledges funding from the Alfred P. Sloan Foundation, as well as from the National Science Foundation under grant SES-1156154.

<sup>†</sup>Ortner: jortner@bu.edu, Chassang: chassang@nyu.edu.

# 1 Introduction

Agents potentially engaging in criminal behavior can undermine institutions by corrupting monitors in charge of investigating them. This paper explores the idea that corruption can be weakened by introducing endogenous asymmetric-information frictions between colluding parties. Building on seminal work by Laffont and Martimort (1997), we show that the cost of deterring crime can be reduced by randomizing the incentives given to the monitor, and letting the magnitude of those incentives serve as the monitor’s private information vis-à-vis the agent. While potential efficiency gains can be significant, the optimality of random incentives depends on pre-existing patterns of asymmetric information. To facilitate policy design, we propose a data-driven framework for prior-free policy evaluation: although aggregate reports by monitors cannot be naïvely used to measure actual criminal activity, we show how to evaluate policy changes using unverified report data. The main takeaway is that a potential local policy change is an improvement if, everything else equal, it is associated with greater reports of crime.

We study a game between three players — a principal, an agent, and a monitor — in which the agent chooses whether or not to engage in criminal behavior  $c \in \{0, 1\}$ . The behavior of the agent is not observed by the principal, but is observed by a monitor who submits report  $m \in \{0, 1\}$ . We think of this report as evidence leading to prosecution: report  $m = 1$  triggers an exogenous judiciary process which imposes a cost  $k$  on criminal agents; report  $m = 0$  (which involves suppression of evidence whenever  $c = 1$ ) triggers no such process. Although the principal cannot observe the agent’s behavior, she can detect misreporting  $m \neq c$  with probability  $q$ . The principal’s only policy control is the efficiency wage  $w$  provided to the monitor.

We allow for collusion between the agent and the monitor at the reporting stage (i.e. corruption). In particular, the monitor can destroy evidence (report message  $m = 0$ ) incriminating a criminal agent in exchange for a bribe. We think of the destruction of evidence as happening in front of the agent, so that there is no moral-hazard between the agent and

the monitor. As a result, collusion boils down to a bilateral trading problem. Exploiting the classic insight that asymmetric information may prevent efficient trade and limits collusion (Myerson and Satterthwaite, 1983, Laffont and Martimort, 1997), we study the extent to which the principal can reduce the cost of incentive provision by creating endogenous asymmetric information between the agent and the monitor.

Our model fits a broad class of environments in which an uninformed principal is concerned about collusion between her monitor and the agents the monitor is supposed to investigate. This includes many of the settings that have been brought up in the empirical literature on corruption, for instance collusion between polluting firms and environmental inspectors (Duflo et al., 2013), tax-evaders and customs officers (Fisman and Wei, 2004), public works contractors and local officials (Olken, 2007), organized crime and police officers (Punch, 2009), and so on. In these settings the principal cannot efficiently monitor agents directly, but may realistically be able to detect tampered evidence by scrutinizing accounts, performing random rechecks in person, or obtaining tips from informed parties. Alternatively, the principal may be able to detect misreporting if crime has delayed but observable consequences, such as environmental pollution, public infrastructure failures, media scandals, and so on.

Our analysis emphasizes three sets of results. The first is that although deterministic incentive schemes are cheap in the absence of collusion, they can become excessively expensive once collusion is allowed. Efficient contracting between the agent and the monitor forces the principal to raise the monitor's wage to the point where the agent and the monitor's joint surplus from misreporting becomes negative. By using random incentives, the principal can reduce the rents of criminal agents, which lowers the cost of incentive provision. We make this point using a simple example without pre-existing asymmetric information. In this case, the cost-savings from using random rather than deterministic incentives are large, in excess of 50% under plausible parameter specifications.

Our second set of results qualifies these optimistic findings by considering environments with pre-existing asymmetric information. In addition to the incentives provided by the

principal, the monitor experiences an exogenous privately-observed idiosyncratic cost  $\eta \geq 0$  for accepting a bribe. We show that the optimality of random incentives depends on the convexity or concavity of the c.d.f.  $F_\eta$  of idiosyncratic costs  $\eta$ . If it is convex over a sufficiently large support, additional asymmetric information is counter-productive.

Finally, motivated by the fact that optimal policy depends crucially on fine details of the environment, we study the possibility of performing prior-free policy evaluations using reporting data from a population of agent-monitor pairs. We consider a principal who has limited knowledge about the parameters of the environment, and hence cannot infer levels of crime from reporting and misreporting data. We first show that aggregate reports of crime across different incentive schemes do not allow for reliable policy evaluation. Indeed, reports of crime depend on both underlying crime rates, and the monitors' decisions to report crime or not. As a result, it is possible that a new incentive scheme decreases aggregate reports of crime, while in fact increasing underlying crime rates. Nevertheless, we show it is possible to perform prior-free local policy evaluations using conditional report data from a single policy (i.e. average reports of crime conditional on realized incentives). Somewhat counter-intuitively, a local policy change improves on a reference incentive scheme if it is associated with higher rates of reported crime. This clarifies that naïvely inferring crime from reporting data leads to incorrect policy recommendations.

This paper and its companion, Chassang and Padró i Miquel (2016), both explore the idea that collusion may be addressed by exploiting informational frictions that make side-contracting difficult. The two papers consider different frictions and emphasize different policy channels. This paper focuses on asymmetric information between the monitor and the agent, and emphasizes endogenous bargaining failures. Chassang and Padró i Miquel (2016) focuses on moral hazard and emphasizes endogenous imperfect monitoring. It departs from the assumption that reports are contractible, so that the monitor is subject to moral hazard. The agent must incentivize her preferred report by committing to a retaliation strategy. To allow information transmission, the principal must limit the information content of her own response to the monitor's reports. Chassang and Padró i Miquel (2016) also attempts to

address the question of policy evaluation. Under the requirement that data from several experiments is available, it shows how to obtain bounds on treatment effects using unverified reports.

On the applied side, this paper relates to and hopes to usefully complement the growing empirical literature on corruption. We address two aspects of the problem which have been emphasized in the literature, for instance in the recent survey by Olken and Pande (2012).<sup>1</sup> The first is that the effectiveness of incentive schemes may be very different over the short-run and the long-run: over time, agents will find ways to corrupt the investigators in charge of monitoring them. We explicitly take into account the possibility of collusion between agents and monitors and propose ways to reduce the cost it imposes on organizations. A second difficulty brought up by Olken and Pande (2012) is that reports of criminal behavior do not provide a reliable measure of underlying crime. Our structural model allows us to back-out measures of underlying crime using observed reports. This connects our work to a small set of papers on structural experiment design (see for instance Karlan and Zinman (2009), Ashraf et al. (2010), Chassang et al. (2012), Chassang and Padró i Miquel (2016), Berry et al. (2012)) that take guidance from structural models to design experiments whose outcome measures can be used to infer unobservable parameters of interest.

On the theory side, our work fits in the literature on collusion in mechanism design initiated by Tirole (1986). It is especially related to Laffont and Martimort (1997, 2000) and Che and Kim (2006, 2009), who emphasize the role of asymmetric information in limiting the extent of collusion.<sup>2</sup> Our contribution is two-fold. First, we show that the principal can potentially benefit from introducing endogenous asymmetric information through random incentives.<sup>3</sup> Second, as a step towards implementation, we show how to evaluate potential

---

<sup>1</sup>For recent work on the measurement of corruption, see Bertrand et al. (2007), and Olken (2007). See also the surveys by Banerjee et al. (2013) and Zitzewitz (2012).

<sup>2</sup>For more on the large literature on collusion in mechanism design, see Felli and Villa-Boas (2000), Faure-Grimaud et al. (2003), Mookherjee and Tsumagari (2004), Burguet and Che (2004), Pavlov (2008), Celik (2009) or Che et al. (2013).

<sup>3</sup>This relates our paper to a recent literature that studies optimal design of information structures; see, for instance, Bergemann and Pesendorfer (2007), Kamenica and Gentzkow (2011), Bergemann et al. (2015), Condorelli and Szentos (2016).

policy changes using only unverified reports. Also related is Baliga and Sjöström (1998), who suggest a distinct mechanism through which random wages (to the agent) may help reduce collusion. They consider a setting in which the agent has no resources of her own, so that any promised payment to the monitor must come from the wage she obtains from the principal. When that is the case, randomizing the agent’s wages undermines her ability to commit to transfers.<sup>4</sup>

Other work has underlined the usefulness of random incentives for reasons unrelated to collusion. In Becker and Stigler (1974) random checks are an optimal response to non-convex monitoring costs. More recently, in work on police crackdowns, Eeckhout et al. (2010) show that in the presence of budget constraints, it may be optimal to provide high powered incentives to a fraction of a population of agents rather than weak incentives to the entire population.<sup>5</sup> In contrast to our analysis, incentives in Eeckhout et al. (2010) must be public information. High powered incentives are useful only if concerned agents are aware of them. In addition, Myerson (1986) and more recently Rahman (2012) emphasize the role of random messaging and random incentives in mechanisms, in particular in settings where the principal needs to disentangle the behavior of different parties.<sup>6</sup>

The paper is organized as follows. Section 2 introduces our framework. Section 3 studies a special case of our model with no pre-existing private information, and delineates the economic forces that make random incentives useful. Section 4 extends the analysis to environments with pre-existing asymmetric information, and shows that additional asymmetric information need not always be optimal. Section 5 proposes an approach to policy-evaluation

---

<sup>4</sup>Also relevant is the work of Basu (2011) and Basu et al. (2014) which highlights the value of asymmetric punishments as a way to make collusion more difficult.

<sup>5</sup>See Lazear (2006) for related results.

<sup>6</sup>Other papers have emphasized the role of random incentives. Rahman and Obara (2010) demonstrate that random messages can improve incentive provision in partnerships by allowing to identify innocent individuals. Jehiel (2012) shows that a principal may benefit from maintaining her agent uninformed about payoff relevant features of the environment, as this may induce higher effort at states at which she values effort most. In a multi-tasking setting, Ederer et al. (2013) show that random contracts may be effective in incentivizing the agent to take a balanced effort profile. In a monopoly pricing context, Calzolari and Pavan (2006a,b) show that a monopolist may benefit from selling to different types of buyers with different probabilities to increase the buyers’ ability to extract revenue on a secondary market.

relying on naturally occurring report data. Section 6 — further developed in the Online Appendix — discusses several extensions to our model including: more sophisticated contracting between the principal and monitor, efficient incomplete-information bargaining between the monitor and agent, extortion from non-criminal agents, and settings in which monitor and agent interact before the agent chooses whether or not to engage in crime. Proofs are collected in Appendix A unless mentioned otherwise.

## 2 Framework

**Players, actions, and payoffs.** We consider a game with three players: a principal, an agent, and a monitor. The agent decides whether to engage in criminal behavior  $c \in \{0, 1\}$ , where crime  $c = 1$  gives the agent a benefit  $\pi_A > 0$  and comes at a cost  $\pi_P < 0$  to the principal. Benefit  $\pi_A > 0$  is the agent’s private information, and is distributed according to a c.d.f.  $F_{\pi_A}$  with density  $f_{\pi_A}$  and support  $[\underline{\pi}_A, \bar{\pi}_A]$ .

The agent’s action is not directly observable to the principal, but is observed by a monitor who chooses to make a report  $m \in \{0, 1\}$  to the principal. We think of this report as evidence leading to prosecution: report  $m = 1$  triggers a judiciary process that imposes an expected cost  $k > \bar{\pi}_A$  on criminal agents and an expected cost  $k_0 \in [0, k]$  on non-criminal agents. This judiciary process is exogenous and outside the control of the principal.

While reports can be falsified (i.e., the monitor can always send either report, regardless of the agent’s action), we assume that the principal detects a false report  $m \neq c$  with probability  $q \in (0, 1)$ , which makes reports partially verifiable. Detection may occur through several channels: for instance accounting discrepancies, random rechecks, or tips from informed parties. Criminal behavior may also have delayed but observable consequences, such as environmental pollution. We further assume that the principal is no longer able to punish a criminal agent after the monitor sends a falsified report  $m = 0$ : the evidence needed for prosecution is no longer available.

The monitor is paid according to a fixed wage contract with wage  $w$ , and gets fired in the

event that the principal finds evidence of misreporting. The monitor is protected by limited liability and cannot be punished beyond the loss of wages.<sup>7</sup> In addition to her expected wage loss, the monitor incurs a cost  $\eta \geq 0$  whenever she misreports. Cost  $\eta$  is the monitor's private information, and is distributed according to a c.d.f.  $F_\eta$  with density  $f_\eta$ .

As part of a possible side-contract, the agent can make transfers  $\tau \geq 0$  to the monitor, i.e. pay her a bribe. Corruption occurs when the monitor accepts to destroy evidence for a criminal agent (i.e. sends message  $m = 0$  although  $c = 1$ ). We assume that crime, rather than corruption, is the behavior that the principal really cares about. Corruption undermines the effectiveness of institutions in charge of punishing crime.

Altogether, expected payoffs  $u_P$ ,  $u_A$ , and  $u_M$  respectively accruing to the principal, the agent, and the monitor take the form:

$$\begin{aligned} u_P &= \pi_P \times c - \gamma_w \times w - \gamma_q \times q \\ u_A &= \pi_A \times c && -[k \times c + k_0 \times (1 - c)] \times m - \tau \\ u_M &= w && -[q \times w + \eta] \times \mathbf{1}_{m \neq c} + \tau, \end{aligned}$$

where  $\gamma_w$  denotes the efficiency cost of raising promised wages and  $\gamma_q$  captures the principal's cost of attention. When the principal is operating under budget or attention constraints, these costs can be interpreted as shadow prices.

We emphasize that the monitor's incentives for truthful reporting are captured by the expected loss from misreporting  $qw + \eta$ . For ease of exposition we treat the distribution of wages  $w$  as the principal's policy variable. However, our analysis applies without change if scrutiny  $q$  is the relevant policy instrument.

**Timing and Commitment.** Our analysis contrasts the effectiveness of incentive schemes under *collusion* and *no-collusion*. The timing of actions is as follows.

1. The principal commits to a distribution of wages  $w$  with c.d.f.  $F_w$ , and draws a random

---

<sup>7</sup>The Online Appendix extends the analysis to the case where the principal and monitor can use arbitrary contracts.



wage  $w$  for the monitor, which is observed by the monitor but not by the agent.

2. The agent chooses whether or not to engage in crime  $c \in \{0, 1\}$ .
3. Under *collusion*, with probability  $\lambda$  the agent makes the monitor a take-it-or-leave-it bribe offer  $\tau$  in exchange for sending message  $m = 0$ ; with probability  $1 - \lambda$  the monitor makes the take-it-or-leave-it bribe offer. We assume perfect commitment so that whenever monitor and agent come to an agreement, the monitor does send message  $m = 0$ . Under *no-collusion* nothing occurs.
4. Under *no-collusion* or, under *collusion* if there was no agreement in the previous stage, the monitor sends the message  $m$  maximizing her final payoff.

We assume for now that parameters  $k$ ,  $k_0$ ,  $\lambda$  and  $q$  are common knowledge. We relax this assumption in Section 5.

We think of non-collusive and collusive environments as respectively capturing short-run and long-run patterns of behavior. In the short run, the agent may take the monitors' behavior as given, and not explore the possibility of bribery. In the long run however, as the agent explores the different options available to her, she may learn that monitors respond favorably to bribes.

**Population interpretation.** Our model admits a natural population interpretation in which distributions  $F_\eta$  and  $F_{\pi_A}$  capture heterogeneity in the population of monitors and agents, and monitors and agents are matched independently. Under this interpretation, wage distribution  $F_w$  captures wage heterogeneity among monitors rather than the randomization of any monitor's wages.

**Motivation.** Our framework is intended to capture the challenges facing public agencies that rely on monitors to assess the behavior of regulated agents. For example, we can think of the principal as an environmental protection agency (EPA), the agent as an industrial plant, and the monitor as an investigator employed by the EPA. In this case, the industrial plant may choose to dump hazardous materials rather than incur the cost of processing

them.<sup>8</sup> Besides environmental protection, other prominent examples include labor safety regulation, tax collection, health inspections, and tackling organized crime. In these cases, crime may respectively correspond to maintaining poor safety and health standards, fraudulent accounting, or extortion and smuggling. The monitor may commit not to report the agent by destroying, or simply by not collecting the evidence needed to initiate a judiciary process. Even if the monitor makes no report of crime, signals of misbehavior may be obtained by the principal after some delay: pollution or poor safety standards may lead to visible consequences (e.g. accidents, local contamination); civil society stakeholders may produce evidence of their own; aggrieved associates of the agent may volunteer incriminating information; and so on . . .

**Modeling assumptions and extensions.** Some of our assumptions are critical to our results, for instance the fact that the principal can commit to a distribution of incentives across monitors, or that monitors cannot verifiably disclose their incentives to agents. We discuss the plausibility of these assumptions in Section 6.

Other assumptions affect the analysis, but do not ultimately change the general thrust of our message. We clarify these assumptions in Section 6 and, when possible, provide appropriate extensions in the Online Appendix. This includes extensions to environments in which bargaining occurs before the agent's crime decision; environments in which the principal can offer the monitor arbitrary contracts; settings in which the monitor and the agent can use arbitrary bargaining mechanisms; as well as environments in which the monitor can extort bribes from non-criminal agents.

---

<sup>8</sup>Note that in the US, environmental pollution is indeed subject to criminal prosecution. The EPA maintains a database of criminal cases resulting from its investigations at <http://www2.epa.gov/enforcement/summary-criminal-prosecutions>.

### 3 Random Wages in a Simple Case

We clarify the potential value of random wages using a simple version of our model in which all monitors have the same cost of falsifying information  $\eta = 0$ , and all agents get the same benefit  $\pi_A < k$  from crime. We further assume that the agent has all the bargaining power at the side-contracting stage, and makes offers with probability  $\lambda = 1$ .

Under these assumptions, the expected cost that a monitor with wage  $w$  incurs from accepting a bribe from a criminal agent and sending a false report is  $qw$ . Thus, under *collusion*, a monitor with wage  $w$  accepts a bribe  $\tau$  from a criminal agent if and only if  $\tau > qw$ .<sup>9</sup> Under *no-collusion*, or if the monitor rejects the agent's offer, the monitor's optimal continuation strategy is to send a truthful report  $m = c$ . In particular, the monitor cannot credibly commit to send report  $m = 1$  when the agent is non-criminal. These observations imply that, under collusion, the expected payoff of a criminal agent of type  $\pi_A$  is  $\pi_A - k + \max_{\tau}(k - \tau)\text{prob}(qw < \tau)$ , and the expected payoff of a non-criminal agent is 0.

**Deterministic wages.** We begin by computing the cost of keeping the agent non-criminal when the principal can use only deterministic wages.

**Lemma 1** (collusion and the cost of incentives). *Assume that the principal uses only deterministic wages. Under no-collusion the principal can induce the agent to be non-criminal at 0 cost.*

*Under collusion, the minimum cost of wages needed to induce the agent to be non-criminal is equal to  $\frac{\pi_A}{q}$ .*

**Proof.** Given any wage  $w$ , under *no-collusion* the monitor's optimal strategy is to send a truthful report. The agent's payoff from action  $c = 1$  is then  $\pi_A - k < 0$  and her payoff from action  $c = 0$  is 0. Thus, under *no-collusion* the principal can induce the agent to be non-criminal at zero cost.

---

<sup>9</sup>By convention, we assume that the monitor rejects the agent's offer whenever she is indifferent between accepting and rejecting a bribe.

Consider next a setting with *collusion*, and note that the monitor accepts a bribe  $\tau$  from a criminal agent if and only if  $\tau > qw$ . The agent's payoff from taking  $c = 1$  is therefore  $\pi_A - \min\{k, qw\}$ , while her payoff from action  $c = 0$  is 0. It follows that the principal can induce the agent to take action  $c = 0$  by setting a deterministic wage  $w = \frac{\pi_A}{q}$ . ■

Lemma 1 shows that, while deterministic incentive schemes work well under no-collusion, their effectiveness is significantly limited whenever collusion is a possibility. We now show that by randomizing wage  $w$  the principal reduces the efficiency of side-contracting between the agent and the monitor, and hence reduces the cost of incentive provision.

**Proposition 1** (optimal incentives under collusion). *Under collusion, the cost-minimizing wage distribution  $F_w^*$  that induces the agent to be non-criminal is described by*

$$\forall w \in [0, \pi_A/q], \quad F_w^*(w) = \frac{k - \pi_A}{k - qw}. \quad (1)$$

The corresponding cost of wages  $W^*(\pi_A) \equiv \mathbb{E}_{F_w^*}[w]$  is

$$W^*(\pi_A) = \frac{\pi_A}{q} \left[ 1 - \frac{k - \pi_A}{\pi_A} \log \left( 1 + \frac{\pi_A}{k - \pi_A} \right) \right] < \frac{\pi_A}{q} \times \frac{\pi_A}{k}. \quad (2)$$

The proof of Proposition 1 is instructive.

**Proof.** A wage distribution  $F_w$  induces the agent to be non-criminal if and only if, for every bribe offer  $\tau \in [0, \pi_A]$ ,  $\pi_A - k + (k - \tau)\text{prob}(\tau > qw) \leq 0$ , or equivalently, if and only if, for every  $\tau \in [0, \pi_A]$ ,  $F_w\left(\frac{\tau}{q}\right) \leq \frac{k - \pi_A}{k - \tau}$ . Using the change in variable  $w = \frac{\tau}{q}$ , we obtain that wage distribution  $F_w$  induces the agent to be non-criminal if and only if,

$$\forall w \in [0, \pi_A/q], \quad F_w(w) \leq \frac{k - \pi_A}{k - qw}. \quad (3)$$

By first-order stochastic dominance, it follows that in order to minimize expected wages, the optimal distribution must satisfy (3) with equality. This implies that the optimal wage

distribution is described by (1). Expected cost expression (2) follows from integration and straightforward computations. ■

Further intuition for why random wages can improve on deterministic wages can be obtained by considering small perturbations around deterministic wage  $\frac{\pi_A}{q}$ . Wage  $\frac{\pi_A}{q}$  deters crime since a criminal agent finds it optimal to offer bribe  $\tau = \pi_A$ , which absorbs all the potential profits from crime. Consider now setting a wage equal to  $\frac{\pi_A}{q}$  with probability  $1 - \epsilon$  and equal to zero otherwise. Since the cost  $k$  of prosecution is strictly higher than  $\pi_A$ , for  $\epsilon > 0$  small enough, a criminal agent will still offer a bribe  $\tau = \pi_A$ . This lets the principal deter crime at a lower expected cost of incentives.

In this simple environment, the savings that can be obtained using random incentives are large: the cost of incentives goes from  $\frac{\pi_A}{q}$  for deterministic mechanisms, to less than  $\frac{\pi_A}{q} \frac{\pi_A}{k}$  for the optimal random incentive scheme. For instance, if the penalty for crime is greater than twice its benefits, i.e.  $k \geq 2\pi_A$ , the principal would be able to save more than 50% on the cost of wages by using random incentives. The gains remain large even if we consider simpler *binary* wage distributions.<sup>10</sup>

**An example.** Binary incentive distributions, boil down to establishing an elite class of harder-to-corrupt monitors. This relates the policies we study to the real-life use of “undercover tactics as a routine part of the inspection process” (Marx, 1992). In one example, *Operation Ampscam*, that took place in New York City, police agents posed as electrical installation inspectors, and arrested contractors who attempted to pay bribes in order to get poor-quality work approved. Undercover police inspectors play the role of hard-to-corrupt monitors in our model. Even if bribing police inspectors is possible, their presence reduces the payoffs of criminal agents by leaving them with two unattractive options. They can either make a high bribe offer that all monitors accept, or make a low bribe offer that undercover

---

<sup>10</sup>For the optimal binary wage distribution, the share of costs saved using random incentives is equal to  $1 - \pi_A/k$ . It puts probability  $1 - \pi_A/k$  on  $w = 0$  and probability  $\pi_A/k$  on  $w = \pi_A/q$ .

police inspectors reject. From the perspective of our model, the fact that Operation Ampscam led to arrests is consistent with the outcome in which criminal agents make low bribe offers that undercover police inspectors reject, and get punished with positive probability.

## 4 Pre-existing Asymmetric Information

**Are random incentives robustly optimal?** The efficiency gains from using random incentives are large in this simple example. Relaxing the assumptions of efficiency wages and take-it-or-leave-it-bargaining does not overturn the optimality of random incentives (see the Online Appendix). Pre-existing asymmetric information poses a more fundamental challenge. Indeed, it is intuitive that complete information should overstate the value of random incentives. Under complete information, random incentives are the only private information allowing the monitor to extract rents from criminal agents.

We return to the general model of Section 2. The monitor experiences a weakly positive private cost  $\eta \sim F_\eta$  for falsifying information, and at the bargaining stage the agent makes the offer with probability  $\lambda$  and the monitor makes the offer with probability  $1 - \lambda$ . Given a distribution of wages  $F_w$ , a criminal agent of type  $\pi_A$  gets an expected payoff equal to

$$U_A(\pi_A) = \pi_A - k + \lambda \max_{\tau \in [0, \pi_A]} (k - \tau) \text{prob}(qw + \eta < \tau).$$

The expression above follows from two observations. First, a monitor with wage  $w$  and type  $\eta$  accepts bribe  $\tau$  from a criminal agent if and only if  $\tau > qw + \eta$ . Second, a monitor demands bribe  $\tau \geq k$  when she acts as proposer at the collusion stage and the agent is criminal, since  $k$  is the highest price criminal agents are willing to pay for a report  $m = 0$ .<sup>11</sup>

As in Section 3, a monitor's optimal continuation strategy is to send a truthful report  $m = c$  if no agreement is reached at the collusion stage. This implies that non-criminal agents get a payoff equal to 0.

---

<sup>11</sup>Specifically, the monitor demands a bribe  $\tau = k$  if  $k \geq qw + \eta$ , and a bribe  $\tau > k$  (which she expects to be rejected) when  $k < qw + \eta$ .

**Policy design under budget constraints.** Given a distribution of wages  $F_w$ , an agent of type  $\pi_A$  will engage in crime if and only if  $U_A(\pi_A) > 0$ . Note that  $U_A(\pi_A)$  is increasing in  $\pi_A$ , so that given a wage profile, agents follow a threshold strategy.

The principal's problem can be decomposed as follows. Given a target threshold  $\pi_A^*$ , find the cheapest wage distribution that implements this threshold. The global optimum can then be found by maximizing over the threshold  $\pi_A^*$ . Alternatively, we can consider the dual problem of a principal who operates under budget constraint  $\mathbb{E}_{F_w}[w] = w_0$ . Given a distribution of wages  $F_w$ , let us denote by  $\bar{\pi}_A(F_w)$  the value of  $\pi_A$  for which an agent is indifferent between actions  $c = 0$  and  $c = 1$ . Given budget  $w_0$ , the principal's problem is to find the distribution of wages  $F_w$  that maximizes threshold  $\bar{\pi}_A(F_w)$  subject to  $\mathbb{E}_{F_w}[w] = w_0$  — this is the *crime-minimizing* wage schedule, given budget  $w_0$ . The overall optimum can then be obtained by optimizing over budget  $w_0$ .

In what follows, we focus on the fixed budget version of the principal's problem. We emphasize that our population interpretation of the model means that the principal can satisfy budget constraint  $\mathbb{E}_{F_w}[w] = w_0$  exactly while using a non-degenerate distribution of wages. The fixed-budget approach is appealing for additional reasons. First, it is realistic: organizations frequently operate within fixed budgets set by other decision-makers. Second, fixed budgets support the principal's ability to commit to mixed strategies. Indeed, taking agent behavior as given, the principal is indifferent over distributions  $\tilde{F}_w$  satisfying  $\mathbb{E}_{\tilde{F}_w}[w] = w_0$ .

**When is additional asymmetric information desirable?** Under pre-existing private information, the optimality of random incentives depends on the shape of distribution  $F_\eta$ .

**Definition 1.** *We say that a wage profile with c.d.f.  $F_w$  is random if and only if the support of  $F_w$  contains at least two elements.*

**Proposition 2** (ambiguous optimal policy). *(i) Whenever  $F_\eta$  is strictly concave over the range  $[0, k]$ , the crime-minimizing wage profile under any budget  $w_0 > 0$  is random.*

(ii) Whenever  $F_\eta$  is strictly convex over the range  $[0, k]$ , the crime-minimizing wage profile under any budget  $w_0 > 0$  is deterministic.

To get some intuition for this result, consider an agent's payoff from taking action  $c = 1$ :

$$\begin{aligned} U_A(\pi_A) &= \pi_A - k + \lambda \max_{\tau \in [0, \pi_A]} (k - \tau) \text{prob}(qw + \eta < \tau) \\ &= \pi_A - k + \lambda \max_{\tau \in [0, \pi_A]} (k - \tau) \mathbb{E}_{F_w}[F_\eta(\tau - qw)]. \end{aligned}$$

If  $F_\eta$  is strictly convex over the support of  $\tau - qw$ , a criminal agent is effectively risk-loving and she obtains a higher payoff from a random wage schedule than from a deterministic one with the same expectation. Inversely, if  $F_\eta$  is strictly concave over the support of  $\tau - qw$ , a criminal agent is effectively risk-averse and her payoff from a random wage schedule is smaller than her payoff from a deterministic one with the same expectation.

If  $F_\eta$  is neither concave nor convex over  $[0, k]$  we can still provide sufficient conditions for random wage profiles to be optimal. Fix a deterministic wage  $w_0 > 0$  and denote by  $\tau_0$  the highest solution to a criminal agent's optimal bribe problem when the monitor is compensated with a deterministic wage  $w_0$ ,

$$\max_{\tau} (k - \tau) \text{prob}(qw_0 + \eta < \tau).$$

**Proposition 3** (sufficient condition for random incentives). *Whenever  $\tau_0 \leq \frac{k}{2}$ , the crime-minimizing policy given budget  $w_0$  is random.*

If starting from a deterministic wage, the agent's optimal bribe is less than half the cost of prosecution, it is optimal to use random wages. The proof exploits the fact that c.d.f.  $F_\eta$  cannot be convex over arbitrarily large ranges of values.<sup>12</sup> The assumption that  $\tau_0 \leq \frac{k}{2}$  lets us exploit non-convexities of  $F_\eta$  around  $w_0$  to construct random wage schedules that improve

---

<sup>12</sup>Note that  $\tau_0 \leq \frac{k}{2}$  implies that  $F_\eta$  is not convex over  $[0, k]$ . Indeed, optimal bribe  $\tau_0$  must satisfy the first-order condition  $f_\eta(\tau_0 - qw_0)(k - \tau_0) = F_\eta(\tau_0 - qw_0)$ . The convexity of  $F_\eta$  over  $[0, k]$  implies that  $F_\eta(\tau_0 - qw_0) \leq f_\eta(\tau_0 - qw_0)(\tau_0 - qw_0) < f_\eta(\tau_0 - qw_0)(k - \tau_0)$ , where the last inequality uses  $\tau_0 \leq \frac{k}{2}$  and  $w_0 > 0$ .



on fixed wages.

We note that when distribution  $F_\eta$  is log-concave, i.e.  $\frac{F_\eta(\cdot)}{f_\eta(\cdot)}$  is increasing, optimal bribe  $\tau_0$  is increasing in  $w_0$ . As a result, the condition of Proposition 3 is more likely to hold when the principal's budget  $w_0$  is small.

Because adding further asymmetric information does not necessarily improve incentive provision, correct policy design must depend on the restrictions, subjective or objective, that the principal can impose on the environment. However, specifying beliefs is often difficult for principals, which makes actual implementation difficult. To address the issue, we show in the next section that it is possible to perform prior-free policy evaluations using naturally occurring unverifiable report data.

## 5 Prior-free Policy Evaluation

We now show that it is possible to evaluate potential local policy changes using reports from monitors under sufficiently rich existing policies. The main takeaway is that marginal policy changes that, everything else equal, increase reports of crime, are local improvements. As a result, naïvely inferring crime from reporting data may lead to incorrect policy recommendations. This result echoes findings from Iyer et al. (2012). In a study of policy changes taking place in India in the early 1990s, the authors show that increased representation of women in local government led to increased reports of crimes against women, but reduced actual crime rates.

**Naïve inference fails.** We first show that a naïve use of reporting data from policy experiments fails to identify the effect that a change in policy has on crime rates.

Consider a principal who is operating under a budget constraint. Given budget  $w_0 > 0$ , let  $F_w^0$  be the deterministic policy under which all monitors are paid wage  $w_0$ , and let  $F_w^1$  be a non-degenerate wage distribution with  $\mathbb{E}_{F_w^1}[w] = w_0$ . We assume that wage policies

$F_w^0$  and  $F_w^1$  are implemented over the same infinite population of exchangeable monitor and agent pairs. We are interested in whether reporting data under the two policies can identify which of them leads to lower crime.

For any policy decision  $d \in \{0, 1\}$ , denote by  $\bar{C}_d$  the proportion of criminal agents under policy  $F_w^d$ . Let  $\bar{R}_d$  be the fraction of monitors reporting  $m = 1$  under policy  $F_w^d$ .<sup>13</sup>

**Lemma 2** (unreliable aggregate reports). *Consider any budget  $w_0 > 0$ , and any random incentive scheme  $F_w^1$  such that  $\mathbb{E}_{F_w^1}[w] = w_0$ .*

*The ordering of reports  $\bar{R}_0$  and  $\bar{R}_1$  is consistent with any ordering of crime  $\bar{C}_0$  and  $\bar{C}_1$ : for any of the four possible pairs of orderings of reports and crime, i.e.,  $\bar{R}_0 \leq \bar{R}_1$  and  $\bar{C}_0 \leq \bar{C}_1$ , there exist specifications of  $k$ ,  $F_{\pi_A}$  and  $F_\eta$  that lead to this ordering.*

In words, the ordering of aggregate reports places no restrictions on the effect of random incentives on crime. Intuitively, reports of crime depend on both the underlying rate of crime and the monitors' decisions to report it. A change in incentive patterns from  $F_w^0$  to  $F_w^1$  changes both the agents' decisions to engage in crime and their bribing behavior. As a result, changes in aggregate reports from  $\bar{R}_0$  to  $\bar{R}_1$  do not always match changes in underlying crime.

**Local policy evaluation.** We now show that an appropriate use of report data from policies with non-degenerate wage distributions can be used to evaluate *local* policy changes. We emphasize three aspects of our results:

- The principal need not to know any of the parameters of the environment: the cost  $k$  imposed by the judiciary on criminal agents, the likelihood  $q$  of detection, and bargaining power  $\lambda$  need not be known.<sup>14</sup>

<sup>13</sup>More explicitly, let  $\bar{\pi}_A(F_w^d)$  denote the type of an agent indifferent between actions  $c = 0$  and  $c = 1$  under policy  $F_w^d$ . Let  $\tau_d$  be a criminal agent's optimal bribe under policy  $F_w^d$ . We have that  $\bar{C}_d = 1 - F_{\pi_A}(\bar{\pi}_A(F_w^d))$  and  $\bar{R}_d = (1 - F_{\pi_A}(\bar{\pi}_A(F_w^d))) \times \text{prob}_{F_w^d}(qw + \eta > \tau_d)$ .

<sup>14</sup>These results contribute to a small literature on mechanism design with limited probabilistic sophistication. This includes maxmin optimal design (Hurwicz and Shapiro, 1978, Hartline and Roughgarden, 2008, Chassang, 2013, Frankel, 2014, Madarász and Prat, 2014, Prat, 2014, Carroll, 2013), as well as data-driven design (Segal, 2003, Chassang and Padró i Miquel, 2016, Brooks, 2014).

- Inference relies on the variation in wages already present in a non-degenerate policy, and does not require knowledge of equilibrium reporting data at the alternative policies.
- If the initial wage distribution is degenerate, a policy experiment is necessary to obtain reporting data for alternative wages. However, it is not necessary to wait for equilibrium crime rates and reports to adjust to the modified wage distribution in order to make policy inferences.

Take as given a non-degenerate wage distribution with cdf  $F_w^0$  and density  $f_w^0$ . We think of distribution  $F_w^0$  as the policy that the principal currently has in place. Let  $f_w^1$  denote a density satisfying

$$\text{supp } f_w^1 \subset \text{supp } f_w^0 \quad \text{and} \quad \mathbb{E}_{f_w^0}[w] = \mathbb{E}_{f_w^1}[w]. \quad (4)$$

When current policy  $f_w^0$  has full support over a range  $[\underline{w}, \bar{w}]$ , the set of policies  $f_w^1$  satisfying (4) is the set of budget-neutral policies with support in  $[\underline{w}, \bar{w}]$ .

For any alternative policy  $f_w^1$  and any  $\epsilon \in [0, 1]$ , construct the mixture  $f_w^\epsilon = (1-\epsilon)f_w^0 + \epsilon f_w^1$ . The proportion of criminal agents under policy  $f_w^\epsilon$  is  $\bar{C}_\epsilon = 1 - F_{\pi_A}(\bar{\pi}_A(f_w^\epsilon))$ , where  $\bar{\pi}_A(f_w^\epsilon)$  is the payoff-type of an agent indifferent between actions  $c = 0$  and  $c = 1$ . We are interested in whether a principal can use reporting data to evaluate the effect that a local policy change in direction  $f_w^1$  (i.e., a marginal increase in  $\epsilon$ ) has on the rate of crime.

Denote by  $\nabla_{f_w^1} \bar{C}$  the gradient of equilibrium crime in policy direction  $f_w^1$ :

$$\nabla_{f_w^1} \bar{C} = \left. \frac{\partial \bar{C}_\epsilon}{\partial \epsilon} \right|_{\epsilon=0}.$$

With this notation, our goal is to evaluate the gradient of crime  $\nabla_{f_w^1} \bar{C}$  for all directions  $f_w^1$ . A marginal move in the direction of  $f_w^1$  is a local policy improvement whenever  $\nabla_{f_w^1} \bar{C} < 0$ .

As an example, suppose the initial policy  $f_w^0$  has support  $\{w_L, w_0, w_H\}$ , with  $w_L < w_0 < w_H$ , and with most of its mass at wage  $w_0$ . Consider a principal who is interested in evaluating whether moving towards a policy with higher variance in incentives will lead to less crime. In this case, policy  $f_w^1$  would be the budget-neutral policy with support  $\{w_L, w_H\}$ .

Let  $\bar{R}_0$  denote the fraction of monitors reporting  $m = 1$  under policy  $f_w^0$ . For any wage  $w \in \text{supp } f_w^0$ , let  $R(w|f_w^0)$  be the fraction of monitors with wage  $w$  reporting  $m = 1$  under the current policy  $f_w^0$ ; i.e.,  $R(w|f_w^0)$  is the share of monitors with wage  $w$  who are matched with a criminal agent and who reject equilibrium bribes under policy  $f_w^0$ .<sup>15</sup> For any policy  $f_w^1$  such that  $\text{supp } f_w^1 \subset \text{supp } f_w^0$  we can construct a counterfactual report of crime under wage distribution  $f_w^1$ , *keeping the agents' behavior constant*, as follows:

$$R_0(f_w^1) \equiv \mathbb{E}_{f_w^0} \left[ R(w|f_w^0) \times \frac{f_w^1(w)}{f_w^0(w)} \right].$$

Counterfactual report  $R_0(f_w^1)$  is the fraction of monitors that would report  $m = 1$  if the principal were to change her policy to  $f_w^1$  and agents continued to behave as if the policy in place was  $f_w^0$ . Counterfactual report  $R_0(f_w^1)$  is obtained by re-weighting reports  $R(w|f_w^0)$  and only requires data from policy  $f_w^0$ . The following result holds.

**Proposition 4** (prior-free policy evaluation). *There exists a fixed coefficient  $\rho > 0$  such that for all alternative policies  $f_w^1$ ,*

$$\nabla_{f_w^1} \bar{C} = \rho [\bar{R}_0 - R_0(f_w^1)].$$

This implies that a small movement from  $f_w^0$  to  $f_w^1$  will decrease crime ( $\nabla_{f_w^1} \bar{C} < 0$ ) if and only if at policy  $f_w^0$ , the counterfactual report of crime reweighted for distribution  $f_w^1$  increases. In other words, it is optimal to move towards the policy  $f_w^1$  such that, everything else equal, would maximize the amount of reported crime. The proof is instructive.

**Proof.** For any policy  $f_w$ , let  $\bar{\pi}_A(f_w)$  be the payoff-type of an agent indifferent between actions  $c = 0$  and  $c = 1$ . Take as given an arbitrary policy  $f_w^1$ . Under wage schedule

---

<sup>15</sup>More explicitly, for all  $w \in \text{supp } f_w^0$ ,  $R(w|f_w^0) = (1 - F_{\bar{\pi}_A}(\bar{\pi}_A(f_w^0))) \times \text{prob}(qw + \eta < \tau_0)$ , where  $\bar{\pi}_A(f_w^0)$  is the cutoff agent type who is indifferent between  $c = 0$  and  $c = 1$  under policy  $f_w^0$ , and  $\tau_0$  is the optimal bribe under policy  $f_w^0$ .

$f_w^\epsilon = (1 - \epsilon)f_w^0 + \epsilon f_w^1$ , the agent's payoff  $U_A^\epsilon(\pi_A)$  from action  $c = 1$  is

$$U_A^\epsilon(\pi_A) = \pi_A - k + \lambda \max_{\tau} (k - \tau) \left[ (1 - \epsilon) \text{prob}_{f_w^0}(qw + \eta < \tau) + \epsilon \text{prob}_{f_w^1}(qw + \eta < \tau) \right].$$

Let  $\tau_0$  be the highest solution to this maximization problem for  $\epsilon = 0$ .

By the Envelope Theorem,  $\forall \pi_A$ ,

$$\begin{aligned} \left. \frac{\partial U_A^\epsilon(\pi_A)}{\partial \epsilon} \right|_{\epsilon=0} &= \lambda(k - \tau_0) \left[ \text{prob}_{f_w^1}(qw + \eta < \tau_0) - \text{prob}_{f_w^0}(qw + \eta < \tau_0) \right] \\ &= \lambda(k - \tau_0) \frac{1}{1 - F_{\pi_A}(\bar{\pi}_A(f_w^0))} \left[ \bar{R}_0 - R_0(f_w^1) \right], \end{aligned}$$

The second equality above follows from two observations. First, mean reports of crime  $\bar{R}_0$  are equal to the product of baseline crime rates times the probability that equilibrium bribes are refused:

$$\bar{R}_0 = [1 - F_{\pi_A}(\bar{\pi}_A(f_w^0))] \times [1 - \text{prob}_{f_w^0}(qw + \eta < \tau_0)].$$

Second, for any  $\tilde{w} \in \text{supp } f_w^0$ , mean reports  $R(\tilde{w}|f_w^0)$  are equal to the product of baseline crime rates times the probability that a monitor with wage  $\tilde{w}$  refuses the equilibrium bribe:

$$\begin{aligned} \forall \tilde{w} \in \text{supp } f_w^0, \quad R(\tilde{w}|f_w^0) &= [1 - F_{\pi_A}(\bar{\pi}_A(f_w^0))] \times [1 - \text{prob}(q\tilde{w} + \eta < \tau_0)] \\ \Rightarrow \quad R_0(f_w^1) &= [1 - F_{\pi_A}(\bar{\pi}_A(f_w^0))] \times [1 - \text{prob}_{f_w^1}(qw + \eta < \tau_0)]. \end{aligned}$$

Since  $\bar{C}_\epsilon = \text{prob}_{F_{\pi_A}}(U^\epsilon(\pi_A) \geq 0) = 1 - F_{\pi_A}(\bar{\pi}_A(f_w^\epsilon))$ , it follows that

$$\begin{aligned} \nabla_{f_w^1} \bar{C} &= \left. \frac{\partial \bar{C}_\epsilon}{\partial \epsilon} \right|_{\epsilon=0} = f_{\pi_A}(\bar{\pi}_A(f_w^0)) \left. \frac{\partial U_A^\epsilon(\pi_A)}{\partial \epsilon} \right|_{\epsilon=0} \\ &= \frac{f_{\pi_A}(\bar{\pi}_A(f_w^0))}{1 - F_{\pi_A}(\bar{\pi}_A(f_w^0))} \lambda(k - \tau_0) \left[ \bar{R}_0 - R_0(f_w^1) \right]. \end{aligned}$$

This proves Proposition 4. ■

**Continuous policy improvement.** The fact that local policy improvements can be identified with naturally occurring data authorizes a process of continuous policy improvement. Starting from a policy  $f_w^0$ , one can engage in gradient-descent by iteratively picking the direction for policy improvement  $f_w^1$  that generates the largest counterfactual report of crime. Note that since accumulated policy-changes cease to be local changes, bribes and crime need to adjust to equilibrium before incremental policy assessments can be made: the process is necessarily gradual.

When  $F_\eta$  is strictly convex over  $[0, k]$  this process pushes initial policy  $f_w^0$  towards fixed deterministic wage  $w_0$ . Indeed, for any policy  $f_w^1$ ,

$$\begin{aligned} \bar{R}_0 - R_0(f_w^1) &= [1 - F_{\pi_A}(\bar{\pi}_A(f_w^0))] \times [\text{prob}_{f_w^1}(qw + \eta < \tau_0) - \text{prob}_{f_w^0}(qw + \eta < \tau_0)] \\ &= [1 - F_{\pi_A}(\bar{\pi}_A(f_w^0))] \times [\mathbb{E}_{f_w^1}[F_\eta(\tau_0 - qw)] - \mathbb{E}_{f_w^0}[F_\eta(\tau_0 - qw)]]. \end{aligned}$$

When  $F_\eta$  is strictly convex over  $[0, k]$ , counterfactual reports  $R_0(f_w^1)$  are maximized by the distribution that puts all its mass point at  $w_0$ . Iterative policy improvement converges to the global policy optimum identified in Proposition 2.

We note that in more general settings, this process (if it converges) will lead to a local policy optimum, rather than a global policy optimum.

**Experiments.** Proposition 4 requires that the support of  $f_w^1$  be included within the support of  $f_w^0$ . When this is not the case, one can obtain an experimental measure  $R_0(f_w^1)$  by randomizing the wage of a small subset of monitors. The proof of Proposition 4 clarifies why one need not wait for equilibrium bribes and crime to adjust in order to interpret the data obtained from such an experiment. Under local policy changes, the equilibrium response of criminal agents has a second order effect on their payoffs. As a result, partial equilibrium responses are sufficient to assess changes in the expected payoffs of crime.

**Evaluating other policy interventions.** The logic of Proposition 4 extends to policy interventions that change truthful-reporting incentives  $qw + \eta$  by affecting the distribution of preference parameter  $\eta$  rather than by changing wages  $w$  or scrutiny  $q$ . For instance, one may consider recruiting monitors from different pools hoping that they may be more or less pro-social. One may also be interested in the effect of a monitor-training, or a morale-enhancing program. In these cases, a local policy change corresponds respectively to marginally increasing the share of monitors recruited from a particular pool, or marginally increasing the share of monitors that undergo the training program. In all these cases, local policy changes towards interventions that yield more reports of crime are policy improvements.

**Caveats.** There are caveats to the policy recommendations following from Proposition 4. The assumptions needed for our results are that: 1) the policy change does not increase the returns  $\pi_A$  to crime; 2) the behavior of a monitor depends only on her *realized* incentives to report truthfully,  $qw + \eta$ .

Hence, Proposition 4 would not be affected if each monitor made an effort decision conditional on her realized incentives  $qw + \eta$ , but it would be affected if the overall policy changed the monitors' propensity to accept bribes. This could happen if monitors as a group found the use of random incentives unfair.<sup>16</sup> Alternatively, policy changes by the principal may cause spite among agents, effectively increasing the returns from crime. This is a concern explored in Iyer et al. (2012), that our framework does not address.

## 6 Discussion

We explored the idea that random incentives can limit the cost of corruption by making side-contracting between criminal agents and monitors more difficult. We show that while the optimality of random incentives depends on unobserved pre-existing patterns of private information, it is possible to use naturally occurring data to guide policy choice. A policy

---

<sup>16</sup>If this is the case, scrutiny  $q$  may be a more appropriate policy variable than wage  $w$ .

change is a local improvement if, everything else equal, it is associated with greater reports of crime. The logic of this result extends to policies that affect truthful-reporting incentives through preferences. Possible implementations of the policies we study are closely related to the use of undercover operations.

The remainder of this section discusses practical aspects of potential implementation as well as alternative modeling choices.

**Commitment and disclosure.** We assume that the principal can commit to a distribution of incentives across the population of monitors, and that monitors cannot disclose their incentives to agents. This is a natural assumption if heterogeneity in truth-telling incentives  $qw$  is created through heterogeneity in scrutiny  $q$ . Attention constraints mean that the principal will focus on a subset of monitors. Furthermore, being under scrutiny is unlikely to be part of a verifiable formal contract.

If wages  $w$  are the relevant policy dimension, commitment to a distribution of wages can be facilitated by first setting an aggregate budget, and then deciding how it should be assigned. This limits the principal's temptation to give all monitors a low wage. In view of the literature on relational contracting (Bull, 1987, Baker et al., 1994, 2002) it is plausible that aspects of compensation, such as promotion or bonuses may not be included in a verifiable contract, but left to the discretion of the principal. Of course greater reliance on the principal's discretion is not without costs, since it creates potential room for abuse on the principal's side.

**Heterogenous incentives without random wages.** The use of heterogeneous wages has distributional implications which stakeholders may find very unfair. This concern can be alleviated while still generating appropriate heterogeneity in incentives.

To the extent that the intensity of scrutiny  $q$  does not affect the welfare of the monitor when she reports truthfully, varying scrutiny  $q$  has limited distributional consequences for non-corrupt monitors. For this reason, it may be a more suitable policy instrument for



practical implementation. Undercover police officers are indeed under much more scrutiny than regular city inspectors. More speculatively, in public infrastructure projects where, as in Olken (2007), local officials play the role of natural monitors, one could vary the probability with which the project gets audited.

Alternatively, one may be able to generate heterogeneous incentives without randomization by letting the monitor's wage depend deterministically on data that is observable to the principal and the monitor, but not the agent. For instance, wages may be contingent on the monitor's tenure, diplomas, the number of crimes she has reported in the past, and so on. Such compensation schemes also introduce heterogeneity in the monitors' incentives, making side-contracting more difficult than under schemes that reward monitors with constant wages.

**Ex ante bargaining.** Our model assumes that the monitor and the agent side-contract after the agent chooses whether to engage in crime. This timing is reasonable in settings where interaction between the monitor and agent is short-lived. For instance, environmental and health inspectors may be rotated across a large number of sites.<sup>17</sup> However, in settings where agents and monitors repeatedly interact, the alternate timing, in which the agent and the monitor bargain before crime happens, may be more plausible.

We show in the Online Appendix that the results of Sections 3 and 4 extend qualitatively under this alternate timing. Endogenous asymmetric information can reduce the costs of incentive provision, but its value depends on pre-existing patterns of asymmetric information. Extending the policy evaluation results of Section 5 is more demanding. The difficulty is that when the monitor and the agent bargain ex ante, there is no report of crime in equilibrium. If they come to an agreement, crime occurs but is not reported. If they do not come to an agreement, crime does not occur. However, we show that reports of attempted corruption, rather than reports of crime, can also be used to evaluate policy. The message is qualitatively the same. Policy changes that increase reports of bribing attempts lower the equilibrium

---

<sup>17</sup>Reasons for rotation, as illustrated by Ohio's EPA 2014 staff rotation initiative, include fostering more homogeneous standards, as well as increasing inspectors' experience.

number of bribing attempts, and reduce the underlying crime rate.<sup>18</sup>

**Extortion.** Our model assumes that the monitor sends a subgame-perfect message following disagreement at the side-contracting stage. This implies that the monitor can never extract bribes from an agent which she observes to be non-criminal. As Olken and Pande (2012) highlight, this prediction is frequently violated: non-criminal agents often have to pay bribes. A simple variation of our baseline model accounts for this. Assume that when the monitor has the bargaining power, she is able to commit to the message she would send in the event of a bargaining failure. A monitor can then extract rents from a non-criminal agent by committing to report the agent as criminal unless a bribe is paid. While this changes the agent’s incentives to engage in crime, we show in the Online Appendix that our main results continue to hold in this setting: random incentives may reduce the cost of corruption, and it is possible to perform local policy evaluation using reporting data.

**Arbitrary contracting between the principal and the monitor.** Throughout the paper we assumed that the monitor is compensated with a fixed wage contract  $w$  and gets fired if she is caught misreporting. Under this assumption, Section 3 shows that deterministic incentive schemes are expensive under collusion, and that the principal can significantly reduce the cost of deterring crime by randomizing the monitor’s wage. These results continue to hold if the principal can use arbitrary contracts to compensate the monitor. With more sophisticated contracts, the principal can reduce the cost of deterring crime by offering the monitor a higher compensation whenever she sends report  $m = 1$ . Indeed, a high compensation following report  $m = 1$  increases the agent’s cost of bribing the monitor, and remains cheap for the principal because it tends to be paid off of the equilibrium path. However, the assumption that reports are only partially verifiable (i.e. false reports are only detected with probability  $q$ ) limits the extent to which the principal can exploit such incentives. With par-

---

<sup>18</sup>Note that the specific results are different. While ranking the prevalence of bribery can be done using bribing-attempts data from a single policy, equilibrium data from two candidate policies is needed to rank crime rates. See the Online Appendix for details.

tially verifiable reports, as the monitor’s compensation following message  $m = 1$  gets large, it becomes optimal for her to report crime regardless of the agent’s action. As a result, the cost of deterring crime with deterministic incentives remains high, and, as we show in the Online Appendix, the cost of keeping the agent non-criminal may be significantly reduced by using random incentives.

**Signaling by the monitor.** One concern with random incentives is that the monitor could signal her type. We address this issue in the Online Appendix by letting the agent and monitor use arbitrary bargaining mechanisms. Because monitors with low-powered incentives benefit from pooling with high-powered monitors, it is impossible for monitors to perfectly signal their types. As a result the principal still benefits from using random incentives.

**Participation constraints.** Throughout the paper we assume that the monitor is risk-neutral, so that randomness in wages does not make participation constraints more difficult to satisfy. Risk-aversion on the monitor’s side may restrain the use of random wages, but our qualitative results continue to hold in that case. The reason for this is that under collusion, participation is not binding. Indeed, in Section 3 we show that the cost of keeping the agent non-criminal with deterministic incentives is equal to  $\frac{\pi_A}{q}$ , compared to an outside option of 0. This means that the principal can use random incentives without affecting the monitor’s participation constraint.

## Appendix

### A Proofs

**Proof of Proposition 2.** The agent’s payoff from taking action  $c = 1$  is

$$\begin{aligned} U_A(\pi_A) &= \pi_A - k + \lambda \max_{\tau \in [0, \pi_A]} (k - \tau) \text{prob}(qw + \eta < \tau) \\ &= \pi_A - k + \lambda \max_{\tau \in [0, \pi_A]} (k - \tau) \mathbb{E}_{F_w} [F_\eta(\tau - qw)]. \end{aligned}$$

Consider first the case in which  $F_\eta$  is strictly concave over  $[0, k]$ . Let  $\tau_0$  be the highest solution to the optimal bribe problem under a deterministic wage  $w_0$  (i.e.,  $\max_\tau (k - \tau)F_\eta(\tau - qw_0)$ ) and note that  $\tau_0 > qw_0$ . Let  $F_w$  be a random wage distribution with  $\mathbb{E}_{F_w}[w] = w_0$  and support  $[w_0 - \gamma, w_0 + \gamma]$ , with  $\gamma > 0$  small enough such that  $\tau_0 > q(w_0 + \gamma)$ . For any  $\epsilon \in [0, 1]$ , let  $F_w^\epsilon = (1 - \epsilon)\mathbf{1}_{w=w_0} + \epsilon F_w$ ; i.e.,  $F_w^\epsilon$  is the mixture between a deterministic wage  $w_0$  and policy  $F_w$ . Since  $F_\eta$  is strictly concave over  $[0, k]$ ,  $(k - \tau)\mathbb{E}_{F_w^\epsilon}[F_\eta(\tau - qw)] < (k - \tau)F_\eta(\tau - qw_0)$  for all  $\tau$  close to  $\tau_0$ . For each  $\epsilon \in [0, 1]$ , let  $\tau_\epsilon$  be the highest solution to  $\max_\tau (k - \tau)\mathbb{E}_{F_w^\epsilon}[F_\eta(\tau - qw)]$ . Since  $\tau_\epsilon$  is close to  $\tau_0$  for  $\epsilon$  small, it follows that

$$(k - \tau_\epsilon)\mathbb{E}_{F_w^\epsilon}[F_\eta(\tau_\epsilon - qw)] < (k - \tau_\epsilon)F_\eta(\tau_\epsilon - qw_0) \leq (k - \tau_0)F_\eta(\tau_0 - qw_0),$$

where the last inequality follows since  $\tau_0$  solves  $\max_\tau (k - \tau)F_\eta(\tau - qw_0)$ . It follows that for  $\epsilon$  small the expected payoff a criminal agent obtains under  $F_w^\epsilon$  is strictly smaller than the one she obtains under the deterministic wage  $w_0$ .

Consider next the case in which  $F_\eta$  is strictly convex over  $[0, k]$ . Note that for any random wage distribution  $F_w$  with  $\mathbb{E}_{F_w}[w] = w_0$ ,  $F_\eta(\cdot)$  is convex over the support of  $\tau - qw$  for all  $\tau \in [0, \pi_A]$ . Therefore, in this case the agent's payoff from being criminal under any random wage distribution with mean  $w_0$  is larger than under the deterministic policy  $w_0$ . ■

**Proof of Proposition 3.** For  $\Delta > 0$ , consider the random wage  $\tilde{w}_\epsilon$  defined by

$$\tilde{w}_\epsilon = \begin{cases} w_0 - \epsilon & \text{with proba } \frac{\Delta}{\Delta + \epsilon} \\ w_0 + \Delta & \text{with proba } \frac{\epsilon}{\Delta + \epsilon}. \end{cases}$$

The expected payoff of a criminal agent under random wage  $\tilde{w}_\epsilon$  is

$$U_A(\pi_A | \tilde{w}_\epsilon) = \pi_A - k + \lambda \max_\tau (k - \tau) \text{prob}_{\tilde{w}_\epsilon}(qw + \eta < \tau).$$

By the Envelope Theorem,

$$\left. \frac{\partial U_A(\pi_A | \tilde{w}_\epsilon)}{\partial \epsilon} \right|_{\epsilon=0} = \lambda(k - \tau_0) \left[ -\frac{1}{\Delta} \text{prob}(qw_0 + \eta < \tau_0) + \frac{1}{\Delta} \text{prob}(q[w_0 + \Delta] + \eta < \tau_0) + qf_\eta(\tau_0 - qw_0) \right].$$

Bribe  $\tau_0$ , which solves  $\max_\tau (k - \tau) \text{prob}(qw_0 + \eta < \tau)$ , must be interior and therefore satisfies the first order condition

$$(k - \tau_0)f_\eta(\tau_0 - qw_0) - \text{prob}(qw_0 + \eta < \tau_0) = 0 \Rightarrow f_\eta(\tau_0 - qw_0) = \frac{\text{prob}(qw_0 + \eta < \tau_0)}{k - \tau_0}.$$

Setting  $\Delta \equiv \tau_0/q - w_0$ , we obtain that

$$\frac{\partial U_A(\pi_A|\tilde{w}_\epsilon)}{\partial \epsilon} \Big|_{\epsilon=0} = q(k - \tau_0)\text{prob}(qw_0 + \eta < \tau_0) \left[ -\frac{1}{\tau_0 - qw_0} + \frac{1}{k - \tau_0} \right] < 0$$

where we used the fact that  $\tau_0 \leq \frac{1}{2}k \Rightarrow k - \tau_0 > \tau_0 - qw_0$ .

Hence for  $\epsilon$  small enough, using random wage distribution  $\tilde{w}_\epsilon$  reduces crime compared to deterministic wage  $w_0$ . ■

**Proof of Lemma 2.** The proof is by example. We proceed case by case and assume throughout that  $\lambda = 1$ . Denote by  $\bar{w}$  and  $\underline{w}$  the maximum and minimum values in the support of  $F_w^1$ . Note that  $w_0 \in (\underline{w}, \bar{w})$ .

We first show that  $\bar{R}_0 < \bar{R}_1$  can be consistent with  $\bar{C}_0 < \bar{C}_1$ . Consider the case where  $k = qw_0$ ,  $F_{\pi_A}$  is a mass point at  $k - \epsilon$  with  $\epsilon > 0$ , and  $F_\eta$  a mass point at 0. For any  $\epsilon > 0$ ,  $\bar{R}_0 = \bar{C}_0 = 0$ . For  $\epsilon > 0$  small enough  $F_w^1(w_0 - \epsilon) > 0$ , which implies that for  $\epsilon$  small enough,

$$\max_{\tau} (k - \tau)\text{prob}_{F_w^1}(qw < \tau) > k - \pi_A = \epsilon.$$

Hence for  $\epsilon > 0$  small enough,  $\bar{C}_1 = 1$ . Furthermore, for  $\epsilon > 0$  small enough,  $F_w^1(w_0 + \epsilon) < 1$ , which implies that  $\bar{R}_1 > 0$  since the agent never offers a bribe  $\tau \geq k = qw_0$ .

Let us show that  $\bar{R}_0 < \bar{R}_1$  can be consistent with  $\bar{C}_0 > \bar{C}_1$ . Set  $F_{\pi_A}$  with full support over  $[0, k]$ , and

$$\eta = \begin{cases} \bar{\eta} & \text{with proba } p \\ 0 & \text{with proba } 1 - p \end{cases}$$

with both  $\bar{\eta} \leq \epsilon$  and  $p \leq \epsilon$ . For  $k$  large enough and  $\epsilon > 0$  small enough, it is immediate that

$$\max_{\tau} (k - \tau)\text{prob}_{F_w^1}(qw + \eta < \tau) < \max_{\tau} (k - \tau)\text{prob}(qw_0 + \eta < \tau)$$

since as  $k$  grows large, it is optimal for the agent to offer bribes respectively converging to  $\bar{w}$  and  $w_0$ , and  $\bar{w} > w_0$ . This implies that  $\bar{C}_0 > \bar{C}_1$ . Let us now show that we can set  $\bar{\eta}$  and  $p$  so that  $\bar{R}_0 < \bar{R}_1$ . A necessary and sufficient condition to obtain  $\bar{R}_0 = 0$  is

$$k - qw_0 - \bar{\eta} > (k - qw_0)(1 - p) \iff k - qw_0 > \frac{\bar{\eta}}{p}. \quad (5)$$

This condition expresses that it is optimal for the agent to offer a bribe  $\tau = qw_0 + \bar{\eta}$  rather than  $\tau = qw_0$  under the deterministic wage  $w_0$ . Similarly, under  $F_w^1$ , a sufficient condition

to ensure that  $\bar{R}_1 > 0$  is that the agent prefer offering a bribe  $\tau = q\bar{w}$  over bribe  $\tau = q\bar{w} + \bar{\eta}$ . A sufficient condition for this is that

$$k - q\bar{w} - \bar{\eta} < (k - q\bar{w})(1 - p) \iff k - q\bar{w} < \frac{\bar{\eta}}{p}. \quad (6)$$

Since  $\bar{w} > w_0$ , it is immediate that for any  $\epsilon$ , one can find values  $p, \bar{\eta} < \epsilon$ , such that conditions (5) and (6) hold simultaneously. For such values,  $\bar{R}_1 > \bar{R}_0 = 0$ , which yields the desired result.

We now show that  $\bar{R}_0 > \bar{R}_1$  can be consistent with  $\bar{C}_0 > \bar{C}_1$ . Set

$$\eta = \begin{cases} \bar{\eta} & \text{with proba } p \\ 0 & \text{with proba } 1 - p \end{cases}$$

with both  $\bar{\eta} \leq \epsilon$  and  $p \leq \epsilon$ . For  $k$  large enough and  $\epsilon > 0$  small enough, we have that

$$\max_{\tau} (k - \tau) \text{prob}_{F_w^1}(qw + \eta < \tau) < \max_{\tau} (k - \tau) \text{prob}(qw_0 + \eta < \tau).$$

Set  $F_{\pi_A}$  as a point mass at a value  $\pi_A$  such that

$$\pi_A - k + \max_{\tau} (k - \tau) \text{prob}_{F_w^1}(qw + \eta < \tau) < 0 < \pi_A - k + \max_{\tau} (k - \tau) \text{prob}(qw_0 + \eta < \tau)$$

for all  $\epsilon$  small enough. This implies that  $\bar{C}_0 = 1 > \bar{C}_1 = 0$ . In turn we obtain that  $\bar{R}_1 = 0$ . Finally, by choosing  $p$  and  $\bar{\eta}$  such that (5) does not hold, one can ensure that  $\bar{R}_0 > 0$ .

Finally, we show that  $\bar{R}_0 > \bar{R}_1$  can be consistent with  $\bar{C}_0 < \bar{C}_1$ . Set  $\eta = 0$ ,  $k = qw_0 - \frac{1}{2}\epsilon$  and

$$\pi_A = \begin{cases} k + \epsilon & \text{with proba } p \\ k & \text{with proba } 1 - p. \end{cases}$$

It is immediate that  $\bar{C}_0 = p$  and  $\bar{R}_0 = p$ . Furthermore, since  $\max_{\tau} (k - \tau) \text{prob}_{F_w^1}(qw + \eta < \tau)$  is strictly positive and bounded away from 0 for  $\epsilon$  small enough, it follows that for  $\epsilon$  small enough  $\bar{C}_1 = 1$  and  $\bar{R}_1 < 1$ . For  $p$  large enough,  $\bar{R}_0 > \bar{R}_1$ . This concludes the proof. ■

## References

ASHRAF, N., J. BERRY, AND J. M. SHAPIRO (2010): “Can Higher Prices Stimulate Product Use? Evidence from a Field Experiment in Zambia,” *The American economic review*,

100, 2383–2413.

BAKER, G., R. GIBBONS, AND K. J. MURPHY (1994): “Subjective Performance Measures in Optimal Incentive Contracts,” *The Quarterly Journal of Economics*, 1125–1156.

——— (2002): “Relational Contracts and the Theory of the Firm,” *Quarterly Journal of Economics*, 39–84.

BALIGA, S. AND T. SJÖSTRÖM (1998): “Decentralization and Collusion,” *Journal of Economic Theory*, 83, 196–232.

BANERJEE, A., S. MULLAINATHAN, AND R. HANNA (2013): *Corruption*, Princeton University Press.

BASU, K. (2011): “Why, for a Class of Bribes, the Act of Giving a Bribe should be Treated as Legal,” .

BASU, K., K. BASU, AND T. CORDELLA (2014): “Asymmetric punishment as an instrument of corruption control,” *World Bank Policy Research Working Paper*.

BECKER, G. S. AND G. J. STIGLER (1974): “Law enforcement, malfeasance, and compensation of enforces,” *J. Legal Stud.*, 3, 1.

BERGEMANN, D., B. BROOKS, AND S. MORRIS (2015): “The Limits of Price Discrimination,” *The American Economic Review*, 105.

BERGEMANN, D. AND M. PESENDORFER (2007): “Information structures in optimal auctions,” *Journal of Economic Theory*, 137, 580–609.

BERRY, J., G. FISCHER, AND R. GUITERAS (2012): “Eliciting and utilizing willingness to pay: evidence from field trials in Northern Ghana,” *Unpublished manuscript*.

BERTRAND, M., S. DJANKOV, R. HANNA, AND S. MULLAINATHAN (2007): “Obtaining a driver’s license in India: an experimental approach to studying corruption,” *The Quarterly Journal of Economics*, 122, 1639–1676.

BROOKS, B. (2014): “Surveying and selling: Belief and surplus extraction in auctions,” .

BULL, C. (1987): “The existence of self-enforcing implicit contracts,” *The Quarterly Journal of Economics*, 147–159.

- BURGUET, R. AND Y.-K. CHE (2004): “Competitive procurement with corruption,” *RAND Journal of Economics*, 50–68.
- CALZOLARI, G. AND A. PAVAN (2006a): “Monopoly with Resale,” *Rand Journal of Economics*, 37, 362–375.
- (2006b): “On the Optimality of Privacy in Sequential Contracting,” *Journal of Economic Theory*, 130, 168–204.
- CARROLL, G. (2013): “Robustness and Linear Contracts,” *Stanford University Working Paper*.
- CELIK, G. (2009): “Mechanism Design with Collusive Supervision,” *Journal of Economic Theory*, 144, 69–75.
- CHASSANG, S. (2013): “Calibrated incentive contracts,” *Econometrica*, 81, 1935–1971.
- CHASSANG, S. AND G. PADRÓ I MIQUEL (2016): “Corruption, Intimidation and Whistleblowing: A Theory of Inference from Unverifiable Reports,” *Unpublished manuscript*.
- CHASSANG, S., G. PADRÓ I MIQUEL, AND E. SNOWBERG (2012): “Selective Trials: A Principal-Agent Approach to Randomized Controlled Experiments,” *American Economic Review*, 102, 1279–1309.
- CHE, Y.-K., D. CONDORELLI, AND J. KIM (2013): “Weak Cartels and Collusion-Proof Auctions,” .
- CHE, Y.-K. AND J. KIM (2006): “Robustly Collusion-Proof Implementation,” *Econometrica*, 74, 1063–1107.
- (2009): “Optimal collusion-proof auctions,” *Journal of Economic Theory*, 144, 565–603.
- CONDORELLI, D. AND B. SZENTES (2016): “Buyer-Optimal Demand and Monopoly Pricing,” Tech. rep., Mimeo.
- DUFLO, E., M. GREENSTONE, R. PANDE, AND N. RYAN (2013): “Truth-telling by Third-party Auditors and the Response of Polluting Firms: Experimental Evidence from India\*,” *The Quarterly Journal of Economics*, 128, 1499–1545.



- EDERER, F., R. HOLDEN, AND M. MEYER (2013): “Gaming and Strategic Ambiguity in Incentive Provision,” *Unpublished manuscript*.
- ECKHOUT, J., N. PERSICO, AND P. E. TODD (2010): “A Theory of Optimal Random Crackdowns,” *American Economic Review*, 100, 1104–1135.
- FAURE-GRIMAUD, A., J.-J. LAFFONT, AND D. MARTIMORT (2003): “Collusion, Delegation and Supervision with Soft Information,” *Review of Economic Studies*, 70, 253–279.
- FELLI, L. AND J. M. VILLA-BOAS (2000): “Renegotiation and Collusion in Organizations,” *Journal of Economics & Management Strategy*, 9, 453–483.
- FISMAN, R. AND S.-J. WEI (2004): “Tax Rates and Tax Evasion: Evidence from ”Missing Imports” in China,” *Journal of Political Economy*, 112.
- FRANKEL, A. (2014): “Aligned delegation,” *The American Economic Review*, 104, 66–83.
- HARTLINE, J. D. AND T. ROUGHGARDEN (2008): “Optimal Mechanism Design and Money Burning,” in *Symposium on Theory Of Computing (STOC)*, 75–84.
- HURWICZ, L. AND L. SHAPIRO (1978): “Incentive structures maximizing residual gain under incomplete information,” *The Bell Journal of Economics*, 9, 180–191.
- IYER, L., A. MANI, P. MISHRA, AND P. TOPALOVA (2012): “The power of political voice: women’s political representation and crime in India,” *American Economic Journal: Applied Economics*, 4, 165–193.
- JEHIEL, P. (2012): “On Transparency in Organizations,” *Unpublished manuscript*.
- KAMENICA, E. AND M. GENTZKOW (2011): “Bayesian Persuasion,” *American Economic Review*, 101, 2590–2615.
- KARLAN, D. AND J. ZINMAN (2009): “Observing unobservables: Identifying information asymmetries with a consumer credit field experiment,” *Econometrica*, 77, 1993–2008.
- LAFFONT, J.-J. AND D. MARTIMORT (1997): “Collusion Under Asymmetric Information,” *Econometrica*, 65, 875–911.
- (2000): “Mechanism Design with Collusion and Correlation,” *Econometrica*, 68, 309–342.

- LAZEAR, E. P. (2006): “Speeding, Terrorism, and Teaching to the Test,” *Quarterly Journal of Economics*, 121, 1029–1061.
- MADARÁSZ, K. AND A. PRAT (2014): “Screening with an Approximate Type Space,” *Working Paper, London School of Economics*.
- MARX, G. T. (1992): “When the guards guard themselves: Undercover tactics turned inward,” *Policing and Society: An International Journal*, 2, 151–172.
- MOOKHERJEE, D. AND M. TSUMAGARI (2004): “The Organization of Supplier Networks: Effects of Delegation and Intermediation,” *Econometrica*, 72.
- MYERSON, R. B. (1986): “Multistage games with communication,” *Econometrica: Journal of the Econometric Society*, 323–358.
- MYERSON, R. B. AND M. A. SATTERTHWAITE (1983): “Efficient mechanisms for bilateral trading,” *Journal of economic theory*, 29, 265–281.
- OLKEN, B. A. (2007): “Monitoring Corruption: Evidence from a Field Experiment in Indonesia,” *Journal of Political Economy*, 115.
- OLKEN, B. A. AND R. PANDE (2012): “Corruption in Developing Countries,” *Annual Review of Economics*, 4, 479–509.
- PAVLOV, G. (2008): “Auction design in the presence of collusion,” *Theoretical Economics*, 3, 383–429.
- PRAT, A. (2014): “Media Power,” *Columbia University Working Paper*.
- PUNCH, M. (2009): *Police corruption: Deviance, accountability and reform in policing*, Routledge.
- RAHMAN, D. (2012): “But Who Will Monitor the Monitor?” *American Economic Review*, 102, 2767–2797.
- RAHMAN, D. AND I. OBARA (2010): “Mediated Partnerships,” *Econometrica*, 78.
- SEGAL, I. (2003): “Optimal pricing mechanisms with unknown demand,” *The American economic review*, 93, 509–529.

TIROLE, J. (1986): “Hierarchies and Bureaucracies: On the Role of Collusion in Organizations,” *Journal of Law, Economics and Organizations*, 2, 181–214.

ZITZEWITZ, E. (2012): “Forensic economics,” *Journal of Economic Literature*, 50, 731–769.

Online Appendix

Making Corruption Harder:  
Asymmetric Information, Collusion, and Crime

Juan Ortner                      Sylvain Chassang\*  
Boston University              New York University

February 28, 2017

**Abstract**

This Online Appendix to “Making Corruption Harder: Asymmetric Information, Collusion, and Crime” provides several extensions. We analyze variants of our baseline model allowing for: ex ante and ex post bargaining, extortion from non-criminal agents, more sophisticated contracting, arbitrary bargaining mechanisms. .

KEYWORDS: monitoring, collusion, corruption, asymmetric information, random incentives, prior-free policy evaluation.

---

\*Ortner: jortner@bu.edu, Chassang: chassang@nyu.edu.

## OA Extensions

### OA.1 Alternative timing of decisions

The model in the main text assumes that the monitor and the agent collude after the agent takes action  $c \in \{0, 1\}$ . This appendix studies the role of random incentives in settings in which the monitor and the agent can collude before the agent chooses her action.

We start by considering a model in which the agent chooses action  $c \in \{0, 1\}$  after side-contracting with the monitor, but which is otherwise the same as the model in Section 3. At the side-contracting stage the agent makes a take-it-or-leave-it offer  $\tau \geq 0$  to the monitor. If the monitor accepts the agent's offer, she commits to send report  $m = 0$  to the principal regardless of the agent's action. Otherwise, if the monitor rejects the agent's offer, she sends the report  $m \in \{0, 1\}$  that maximizes her expected payoff. The principal detects false messages with probability  $q$ . The monitor is compensated with an efficiency wage  $w \geq 0$ , and loses this wage if the principal detects that the message was false. We assume for now that all monitors have a type  $\eta = 0$  and that all agents have type  $\pi_A < k$ . We relax these assumptions later.

**Lemma OA.1.** *The agent takes action  $c = 1$  if and only if the monitor accepts her bribe.*

**Proof.** If the monitor accepts the agent's bribe  $\tau$ , the agent's payoffs from action  $c = 1$  is  $\pi_A - \tau$ , while her payoff from action  $c = 0$  is  $-\tau$ . If the monitor rejects the agent's bribe, the agent's payoff from  $c = 1$  is  $\pi_A - k < 0$  (since in this case the monitor will find it optimal to send message  $m = 1$ ), while her payoff from action  $c = 0$  is 0. Therefore, the agent takes action  $c = 1$  if and only if the monitor accepts her bribe. ■

Lemma OA.1 implies that the monitor's payoff from accepting bribe  $\tau$  is  $\tau + (1 - q)w$ , while her payoff from rejecting the bribe and sending a truthful message is  $w$ . Therefore, a monitor with wage  $w$  accepts bribe  $\tau$  if and only if  $\tau > qw$ .

We now consider the case in which the principal compensates the agent with a determin-

istic wage  $w$ . The following result generalizes Lemma 1 to the current setting; its proof is identical to the proof of Lemma 1 and hence omitted.

**Lemma OA.2.** *Suppose the principal uses a deterministic wage  $w$ . Under collusion, the minimum cost of wages needed to induce the agent to take action  $c = 0$  is equal to  $\frac{\pi_A}{q}$ .*

Consider next the case in which the principal randomizes over the monitor's wage. Suppose the principal pays the monitor an efficiency wage drawn from the c.d.f.  $F$ . Note that the agent's payoff from making an offer  $\tau \geq 0$  is  $F(\tau/q) \times (\pi_A - \tau) + (1 - F(\tau/q)) \times 0$ . Let  $\tau_F^*$  be the smallest solution to  $\max_{\tau} F(\tau/q)(\pi_A - \tau)$ . For any distribution  $F$ , the principal's expected payoff is

$$F\left(\frac{\tau_F^*}{q}\right) \pi_P - \gamma_w \mathbb{E}_F[w] - \gamma_q q.$$

Under wage distribution  $F$ , the monitor accepts the agent's bribe when her wage is lower than  $\tau_F^*/q$ . In this case, the agent takes action  $c = 1$  and the principal incurs cost  $\pi_P < 0$ .

**Proposition OA.1.** *Assume that the agent and monitor collude before the agent chooses  $c \in \{0, 1\}$ . Then, the optimal wage distribution  $\tilde{F}^*$  is described by,*

$$\forall w \in \left[0, \frac{\pi_A}{q} \left(1 - e^{-\frac{q}{\gamma_w} \frac{\pi_P}{\pi_A}}\right)\right], \quad \tilde{F}_w^*(w) = \frac{e^{\frac{q}{\gamma_w} \frac{\pi_P}{\pi_A}} \pi_A}{\pi_A - qw}. \quad (\text{O1})$$

*When the principal pays the monitor a wage drawn from  $\tilde{F}_w^*$ , the agent takes action  $c = 1$  with probability  $\tilde{F}_w^*(0) \in (0, 1)$ .*

**Proof.** Consider first distributions  $F$  such that  $F\left(\frac{\tau_F^*}{q}\right) = 0$ . Note that  $F\left(\frac{\tau_F^*}{q}\right) = 0$  implies that  $0 \geq \max_{\tau} F(\tau/q)(\pi_A - \tau)$ , and so  $F(\tau/q) = 0$  for all  $\tau < \pi_A$ . Therefore, for distributions  $F$  such that  $F\left(\frac{\tau_F^*}{q}\right) = 0$ , the minimum cost of wages is achieved with a distribution that puts all its mass at  $w = \pi_A/q$ . The principal's payoff under this distribution is  $-\gamma_w \frac{\pi_A}{q} - \gamma_q q$ . Our arguments below show that such a distribution is never optimal.

Consider next distributions  $F$  such that  $F\left(\frac{\tau_F^*}{q}\right) > 0$ . Since  $\tau_F^* \geq 0$  is the optimal offer,

for all  $\tau \geq 0$ ,

$$F\left(\frac{\tau_F^*}{q}\right)(\pi_A - \tau_F^*) \geq F\left(\frac{\tau}{q}\right)(\pi_A - \tau) \iff F\left(\frac{\tau}{q}\right) \leq F\left(\frac{\tau_F^*}{q}\right) \frac{\pi_A - \tau_F^*}{\pi_A - \tau}. \quad (\text{O2})$$

By first order stochastic dominance, an optimal wage distribution  $F$  with  $F\left(\frac{\tau_F^*}{q}\right) > 0$  must be such that (O2) holds with equality for all  $\tau$  such that  $F(\tau/q) < 1$ .

Next, we show that the optimal distribution  $F$  with  $F\left(\frac{\tau_F^*}{q}\right) > 0$  must be such that  $\tau_F^* = 0$ . Let  $F$  be such that  $\tau_F^* > 0$ , and let  $\hat{F}$  be an alternative distribution described by:  $\hat{F}(0) = F(\tau_F^*/q)$  and  $\hat{F}(\tau/q) = \frac{\hat{F}(0)\pi_A}{\pi_A - \tau}$  for all  $\tau \in [0, \pi_A(1 - \hat{F}(0))]$ . By construction, bribe  $\tau = 0$  maximizes  $\hat{F}(\tau/q)(\pi_A - \tau)$ . Since  $\hat{F}(0) = F(\tau_F^*/q)$ , the probability that the agent takes action  $c = 1$  is the same under  $\hat{F}$  than under  $F$ . Moreover, for all  $\tau$  such that  $\hat{F}(\tau/q) < 1$ ,  $\hat{F}(\tau/q) = \hat{F}(0) \frac{\pi_A}{\pi_A - \tau} > F(\tau_F^*/q) \frac{\pi_A - \tau_F^*}{\pi_A - \tau} \geq F(\tau/q)$  (where the last inequality follows since offer  $\tau_F^*$  is optimal under policy  $F$ ). This implies that  $\mathbb{E}_F[w] > \mathbb{E}_{\hat{F}}[w]$ , so the principal's payoff is larger under  $\hat{F}$  than under  $F$ .

Using the change in variable  $w = \tau/q$ , the two paragraphs above imply that the optimal wage distribution  $F$  with  $F\left(\frac{\tau_F^*}{q}\right) > 0$  is such that  $\tau_F^* = 0$  and is described by

$$\forall w \in \left[0, \frac{\pi_A}{q}(1 - F(0))\right], \quad F(w) = \frac{F(0)\pi_A}{\pi_A - qw}.$$

The principal's expected payoff from using this wage distribution is

$$F(0)\pi_P - \gamma_w \mathbb{E}_F[w] - \gamma_q q = F(0)\pi_P - \gamma_w \frac{\pi_A}{q}(1 - F(0) + F(0) \ln F(0)) - \gamma_q q.$$

This expression is strictly concave in  $F(0)$ , and converges to  $-\gamma_w \frac{\pi_A}{q} - \gamma_q q$  as  $F(0) \rightarrow 0$ . Maximizing this expression with respect to  $F(0)$  yields  $F(0) = e^{\frac{q}{\gamma_w} \frac{\pi_P}{\pi_A}} \in (0, 1)$ . Therefore, the optimal wage distribution is given by (O1). ■

Proposition OA.1 shows that random incentives are optimal in this setting. We note

that the principal can improve upon deterministic wages using simpler schemes. Suppose the principal uses a two-wage distribution, paying the monitor wage  $\underline{w} = 0$  with probability  $x \in [0, 1]$  and wage  $\bar{w} = \frac{\pi_A}{q}(1 - x)$  with probability  $1 - x$ . Under this wage distribution, it is optimal for the agent to make a bribe offer of  $\tau = 0$ . The principal's payoff under this distribution is  $x\pi_p - \gamma_w(1 - x)^2 \frac{\pi_A}{q} - \gamma_q q$ , which is maximized by setting  $x = \max\{0, 1 + \frac{q}{\gamma_w} \frac{\pi_p}{2\pi_A}\}$ .

**Ambiguous optimal policy.** Next, we extend Proposition 2 to this environment. As in Section 4, we assume that monitors and agents are privately informed about their types, with  $\eta$  distributed according to c.d.f.  $F_\eta$  with density  $f_\eta$  and  $\pi_A$  distributed according to c.d.f.  $F_{\pi_A}$  with density  $f_{\pi_A}$ .

Given wage distribution  $F_w$ , an agent with type  $\pi_A$  offers bribe  $\tau$  solving

$$\begin{aligned} U(\pi_A) &= \max_{\tau \in [0, \pi_A]} \text{prob}_{F_w}(qw + \eta < \tau)(\pi_A - \tau) \\ &= \max_{\tau \in [0, \pi_A]} \mathbb{E}_{F_w}[F_\eta(\tau - qw)](\pi_A - \tau). \end{aligned} \quad (\text{O3})$$

Equation (O3) can be used to extend Proposition 2 to this environment. Indeed, whenever  $F_\eta$  is strictly concave (strictly convex) over the range  $[0, \pi_A]$ , the wage profile that minimizes the agent's payoff under any budget  $w_0 > 0$  is random (deterministic). Note, however, that these statements relate to the agent's payoff, and not to the probability that the agent is criminal. It is also possible to find conditions on  $F_\eta$  under which the crime-minimizing policy is deterministic. For instance, if  $F_\eta$  and  $f_\eta$  are both strictly convex, and  $f_\eta(\tau - qw_0)(\pi_A - \tau)$  is strictly decreasing in  $\tau$ , then the crime-minimizing policy is deterministic.<sup>1</sup>

---

<sup>1</sup>Proof: Fix a budget  $w_0$  and let  $F_w$  be any random policy with  $\mathbb{E}_{F_w}[w] = w_0$ . Let  $\tau_0$  be the highest solution to  $\max_\tau (\pi_A - \tau)F_\eta(\tau - qw_0)$  and  $\tau_{F_w}$  be the highest solution to  $\max_\tau \mathbb{E}_{F_w}[F_\eta(\tau - qw)](\pi_A - \tau)$ . Suppose by contradiction that the probability with which the agent is criminal is higher under the deterministic policy than under policy  $F_w$ , so  $F_\eta(\tau_0 - qw_0) \geq \mathbb{E}_{F_w}[F_\eta(\tau_{F_w} - qw)]$ . Note that  $F_\eta$  strictly convex implies  $\tau_0 > \tau_{F_w}$ . Then,

$$(\pi_A - \tau_0)f_\eta(\tau_0 - qw_0) = F_\eta(\tau_0 - qw_0) \geq \mathbb{E}_{F_w}[F_\eta(\tau_{F_w} - qw)] = (\pi_A - \tau_{F_w})\mathbb{E}_{F_w}[f_\eta(\tau_{F_w} - qw)],$$

where the first and last equalities follow since  $\tau_0$  and  $\tau_{F_w}$  are optimal and satisfy the first-order conditions. Finally, since  $\tau_0 > \tau_{F_w}$  and  $f_\eta(\tau - qw_0)(\pi_A - \tau)$  is strictly decreasing in  $\tau$ , the inequality above implies that  $(\pi_A - \tau_{F_w})f_\eta(\tau_{F_w} - qw_0) > (\pi_A - \tau_{F_w})\mathbb{E}_{F_w}[f_\eta(\tau_{F_w} - qw)]$ , which cannot be since  $f_\eta$  is strictly convex. Hence,



**Policy evaluation.** We now show how the policy evaluation results in Section 5 extend to an environment in which the interaction between monitors and agents may be ex-ante or ex-post. In particular, we consider a model in which a fraction  $\mu \in (0, 1)$  of agents interact with their monitors after taking action  $c \in \{0, 1\}$ , as in the main text, and a fraction  $1 - \mu$  of agents interact with their monitors before taking action  $c \in \{0, 1\}$ . Fraction  $\mu$  is unknown to the principal. We assume that the agent has all the bargaining power at the side-contracting stage, and makes offers with probability 1.<sup>2</sup>

We allow monitors to report failed bribing attempts, in addition to reports of crime: monitors now send crime reports  $m_c \in \{0, 1\}$  and bribing attempt reports  $m_b \in \{0, 1\}$ . As in the baseline model, an agent who was reported  $m_c = 1$  incurs a cost of  $k$  if criminal, and a cost  $k_0 \leq k$  if not criminal. In addition, an agent who was reported  $m_b = 1$  incurs a small fine  $\phi > 0$  if she was not reported for crime; if she was reported  $m_c = 1$ , she incurs cost  $k$  if criminal and cost  $k_0$  if not criminal. We note that allowing monitors to report bribing attempts is needed to generate variation in the monitors' reports that can be used to evaluate how different policies affect crime among those agents that interact with monitors ex-ante. Indeed, by Lemma OA.1, agents who side-contract with monitors ex-ante take action  $c = 1$  if and only if their monitor accepts the bribe. As a result, monitors who interact ex-ante with agents always report  $m_c = 0$  regardless of the policy in place.

We start by considering agents who interact with monitors ex-ante. Given wage distribution  $F_w$ , the expected payoff of an agent with type  $\pi_A$  who interacts with her monitor ex-ante and who engages in bribing behavior is

$$\begin{aligned} U_{F_w}^{\text{ante}}(\pi_A) &= \max_{\tau \in [0, \pi_A]} \text{prob}_{F_w}(qw + \eta < \tau)(\pi_A - \tau) + \text{prob}_{F_w}(qw + \eta > \tau)(-\phi) \\ &= \max_{\tau \in [0, \pi_A]} \text{prob}_{F_w}(qw + \eta < \tau)(\pi_A + \phi - \tau) - \phi. \end{aligned}$$

Such an agent will engage in bribing behavior if and only if  $U_{F_w}^{\text{ante}}(\pi_A) > 0$ ; if she engages

---

the crime-minimizing policy is deterministic.

<sup>2</sup>Our results extend to a setting with probabilistic take-it-or-leave-it offers provided that the monitor observes the agent's type.

in bribing behavior, she takes action  $c = 1$  if and only if her bribe is accepted. We note that monitors with type  $\eta > 0$  who interact with an agent ex-ante have a strict incentive to report failed bribing attempts. As a result, with probability 1, a monitor who interacts ex-ante with an agent engaging in bribing behavior will report  $m_b = 0$  if she accepts the agent's bribe, and  $m_b = 1$  if she rejects it. By our arguments above, monitors who interact ex-ante always report  $m_c = 0$ .

Consider next agents who interact with monitors ex-post. Given policy  $F_w$ , the expected payoff of a criminal agent of type  $\pi_A$  who interacts with her monitor ex-post is

$$U_{F_w}^{\text{post}}(\pi_A) = \pi_A - k + \max_{\tau \in [0, k]} (k - \tau) \text{prob}_{F_w}(qw + \eta < \tau).$$

An agent of type  $\pi_A$  who interacts with her monitor ex-post chooses  $c = 1$  if and only if  $U_{F_w}^{\text{ante}}(\pi_A) > 0$ . A monitor who interacts with a criminal agent ex-post reports  $m_c = m_b = 0$  if she accepts the bribe, and reports  $m_c = m_b = 1$  if she rejects it (by assumption, in the latter case the agent incurs a punishment cost of  $k$ ). On the other hand, a monitor who interacts with a non-criminal agent reports  $m_b = m_c = 0$ .

We now show how a principal can use reports from failed bribing attempts to perform local policy evaluations on agents who interact with monitors ex-ante. Take as given a wage distribution with c.d.f.  $F_w^0$  and density  $f_w^0$ , and let  $f_w^1$  be a policy with  $\text{supp } f_w^1 \subset \text{supp } f_w^0$  and  $\mathbb{E}_{f_w^0}[w] = \mathbb{E}_{f_w^1}[w]$ . For any such policy  $f_w^1$  and any  $\epsilon \in [0, 1]$ , construct the mixture  $f_w^\epsilon = (1 - \epsilon)f_w^0 + \epsilon f_w^1$ .

Given policy  $f_w^\epsilon$ , we denote by  $\bar{R}_\epsilon^b(\pi_A)$  the proportion of monitors who report  $m_b = 1$  and  $m_c = 0$  among monitors matched with an agent of type  $\pi_A$ . We denote by  $U_{f_w^\epsilon}^{\text{ante}}(\pi_A)$  the payoff of an agent of type  $\pi_A$  from engaging in bribing behavior:

$$U_{f_w^\epsilon}^{\text{ante}}(\pi_A) = \max_{\tau} \text{prob}_{f_w^\epsilon}(qw + \eta < \tau)(\pi_A + \phi - \tau) - \phi$$

Fix a type  $\pi_A$  such that  $U_{f_w^0}^{\text{ante}}(\pi_A) > 0$ . For any  $f_w^1$ , denote by  $\nabla_{f_w^1} U(\pi_A)$  the gradient of the

agent's payoff from bribing in policy direction  $f_w^1$ :

$$\nabla_{f_w^1} U^{\text{ante}}(\pi_A) = \left. \frac{\partial U_{f_w^\epsilon}^{\text{ante}}(\pi_A)}{\partial \epsilon} \right|_{\epsilon=0}.$$

For any wage  $w \in \text{supp } f_w^0$ , let  $R_0^b(w; \pi_A)$  be the fraction of reports  $m_b = 1$  and  $m_c = 0$  from monitors with wage  $w$  who were matched with agents of type  $\pi_A$  under policy  $f_w^0$ . For any  $f_w^1$  with  $\text{supp } f_w^1 \subset \text{supp } f_w^0$ , construct counterfactual reports

$$R_0^b(f_w^1; \pi_A) \equiv \mathbb{E}_{f_w^0} \left[ R_0^b(w; \pi_A) \times \frac{f_w^1(w)}{f_w^0(w)} \right]. \quad (\text{O4})$$

The following result holds.

**Proposition OA.2.** *For every  $\pi_A$  with  $U_{f_w^0}^{\text{ante}}(\pi_A) > 0$ , there exists a fixed coefficient  $\rho(\pi_A) > 0$  such that for all alternative policies  $f_w^1$ ,*

$$\nabla_{f_w^1} U(\pi_A) = \rho(\pi_A) \left[ \bar{R}_0^b(\pi_A) - R_0^b(f_w^1; \pi_A) \right].$$

**Proof.** Take as given a policy  $f_w^1$ . Under wage schedule  $f_w^\epsilon$ , the payoff of an agent with type  $\pi_A$  who engages in bribing behavior is

$$U_{f_w^\epsilon}^{\text{ante}}(\pi_A) = \max_{\tau} (\pi_A + \phi - \tau) [(1 - \epsilon) \text{prob}_{f_w^0}(qw + \eta < \tau) + \epsilon \text{prob}_{f_w^1}(qw + \eta < \tau)] - \phi.$$

Let  $\tau_0$  be the highest solution to this maximization problem for  $\epsilon = 0$ . By the Envelope Theorem,

$$\begin{aligned} \nabla_{f_w^1} U(\pi_A) &= (\pi_A + \phi - \tau_0) [\text{prob}_{f_w^1}(qw + \eta < \tau_0) - \text{prob}_{f_w^0}(qw + \eta < \tau_0)] \\ &= (\pi_A + \phi - \tau_0) \frac{1}{1 - \mu} \left[ \bar{R}_0^b(\pi_A) - R_0^b(f_w^1; \pi_A) \right], \end{aligned} \quad (\text{O5})$$

where  $1 - \mu \in (0, 1)$  is the fraction of agents that interact with monitors ex-ante. The second equality above follows from two observations. First, mean reports of failed bribing

attempts  $\bar{R}_0^b(\pi_A)$  are equal to the product of the fraction of agents of type  $\pi_A$  who interact with monitors ex-ante times the probability that their equilibrium bribes are refused:

$$\bar{R}_0^b(\pi_A) = (1 - \mu) \times [1 - \text{prob}_{f_w^0}(qw + \eta < \tau_0)].$$

Second, for any  $\tilde{w} \in \text{supp } f_w^0$ , mean reports  $R_0^b(w; \pi_A)$  are equal to the product of the fraction of agents of type  $\pi_A$  who interact with monitors ex-ante times the probability that a monitor with wage  $\tilde{w}$  refuses their bribe:

$$\begin{aligned} \forall \tilde{w} \in \text{supp } f_w^0, \quad R_0^b(w; \pi_A) &= (1 - \mu) \times [1 - \text{prob}(q\tilde{w} + \eta < \tau_0)] \\ \Rightarrow \quad R_0^b(f_w^1, \pi_A) &= (1 - \mu) \times [1 - \text{prob}_{f_w^1}(qw + \eta < \tau_0)]. \end{aligned}$$

This establishes the result. ■

Proposition OA.2 shows that, under this alternative timing, a principal who can condition on the type of the agent can evaluate how small changes in policy affect the agent's payoff from engaging in bribing behavior. We note that, even when the agent's type is unobservable, the identification result in Proposition OA.2 can still be useful if the principal can condition on a sufficiently rich set of covariates.

Proposition OA.2 can be used to identify directions of policy change that lead to less bribing behavior. We now show how this result can be leveraged to evaluate the effect of local policy changes on crime rates among agents who interact ex-ante.

Let  $f_w^0$  be the original policy in place. For any policy  $f_w$ , we let  $\bar{\pi}_A^{\text{ante}}(f_w)$  denote the cutoff such that all agents with  $\pi_A > \bar{\pi}_A^{\text{ante}}(f_w)$  who interact ex-ante engage in bribing behavior under policy  $f_w$ . Let  $f_w^1$  be a policy direction that reduces the set of agents who engage in bribing behavior; i.e., a policy direction with  $\nabla_{f_w^1} U(\bar{\pi}_A^{\text{ante}}(f_w^0)) < 0$ . Fix  $\epsilon > 0$  small, and let  $f_w^\epsilon = (1 - \epsilon)f_w^0 + \epsilon f_w^1$ . Suppose that the principal changes her policy from  $f_w^0$  to  $f_w^\epsilon$ .

Let  $\bar{C}_0^{\text{ante}}$  and  $\bar{C}_\epsilon^{\text{ante}}$  denote, respectively, the fraction of agents who interact ex-ante that

take action  $c = 1$  under policies  $f_w^0$  and  $f_w^\epsilon$ . Let  $\bar{R}_0^b$  and  $\bar{R}_\epsilon^b$  denote, respectively, the fraction of monitors who report bribing attempts and don't report crime under policies  $f_w^0$  and  $f_w^\epsilon$ ; i.e., the fraction of monitors who report  $m_b = 1$  and  $m_c = 0$ . The following result holds.

**Proposition OA.3.** *Fix a policy  $f_w^0$  and a policy direction  $f_w^1$  such that  $\nabla_{f_w^1} U(\bar{\pi}_A^{\text{ante}}(f_w^0)) < 0$ . Then, there exists a constant  $\kappa > 0$  such that*

$$\bar{C}_\epsilon^{\text{ante}} - \bar{C}_0^{\text{ante}} \leq \kappa \times [\bar{R}_0^b - \bar{R}_\epsilon^b].$$

**Proof.** For every type  $\pi_A$  and any policy  $f_w$ , let  $\tau(\pi_A; f_w)$  denote the bribe that agents of type  $\pi_A > \bar{\pi}_A^{\text{ante}}(f_w)$  who interact with monitors ex-ante offer under policy  $f_w$ ; i.e.,  $\tau(\pi_A; f_w)$  maximizes  $\text{prob}_{f_w}(qw + \eta < \tau)(\pi_A + \phi - \tau)$ . Note that an agent of type  $\pi_A > \bar{\pi}_A^{\text{ante}}(f_w)$  who interacts ex-ante takes action  $c = 1$  only when her bribe is accepted. Then, for  $x \in \{0, \epsilon\}$ , the fraction of agents who interact ex-ante that take action  $c = 1$  under policy  $f_w^x$  is

$$\bar{C}_x^{\text{ante}} = (1 - F_{\pi_A}(\bar{\pi}_A^{\text{ante}}(f_w^x))) \times \mathbb{E}_{F_{\pi_A}}[\text{prob}_{f_w^x}(qw + \eta < \tau(\pi_A; f_w^x)) | \pi_A > \bar{\pi}_A^{\text{ante}}(f_w^x)]. \quad (\text{O6})$$

Note next that, under policy  $f_w$ , an agent of type  $\pi_A$  who interacts ex-ante gets reported for a failed bribing attempt with probability  $\mathbf{1}_{\{\pi_A > \bar{\pi}_A^{\text{ante}}(f_w)\}} \times (1 - \text{prob}_{f_w}(qw + \eta < \tau(\pi_A; f_w)))$ . Moreover, only monitors who interact with agents ex-ante send reports  $m_b = 1$  and  $m_c = 0$ .<sup>3</sup> Thus, for  $x \in \{0, \epsilon\}$ , the share of monitors reporting  $m_b = 1$  and  $m_c = 0$  under policy  $f_w^x$  is

$$\begin{aligned} \bar{R}_x^b &= (1 - \mu) \times (1 - F_{\pi_A}(\bar{\pi}_A^{\text{ante}}(f_w^x))) \times \mathbb{E}_{F_{\pi_A}}[1 - \text{prob}_{f_w^x}(qw + \eta < \tau(\pi_A; f_w^x)) | \pi_A > \bar{\pi}_A^{\text{ante}}(f_w^x)] \\ &= (1 - \mu) \times [(1 - F_{\pi_A}(\bar{\pi}_A^{\text{ante}}(f_w^x))) - \bar{C}_x^{\text{ante}}], \end{aligned} \quad (\text{O7})$$

where we used equation (O6). Since policy direction  $f_w^1$  is such that  $\nabla_{f_w^1} U(\bar{\pi}_A^{\text{ante}}(f_w^0)) < 0$ , it

---

<sup>3</sup>Monitors who interact with agents ex-post send either  $m_b = m_c = 0$  or  $m_b = m_c = 1$ .

follows that  $\bar{\pi}_A^{\text{ante}}(f_w^0) < \bar{\pi}_A^{\text{ante}}(f_w^\epsilon)$ . Using this together with equation (O7) yields

$$\bar{C}_\epsilon^{\text{ante}} - \bar{C}_0^{\text{ante}} \leq \frac{1}{1-\mu} [\bar{R}_0^b - \bar{R}_\epsilon^b],$$

which establishes the result. ■

We end this section by noting that the principal can perform local policy evaluations on agents who interact ex-post using reports of crime. Indeed, since reports  $m_c = 1$  come exclusively from monitors who interact with agents ex-post, Proposition 4 continues to hold in this setting. This, combined with Propositions OA.2 and OA.3, allows the principal to find policy directions that reduce overall crime rates.

## OA.2 Extortion

This section shows how our results extend to settings in which the monitor can extort transfers from non-criminal agents by committing to send a false report. The framework we consider is essentially the same as in Section 4. The only difference is that a monitor who makes an offer at the side-contracting stage can commit to sending a false report if the agent rejects her proposal. A report  $m = 1$  triggers an exogenous judiciary process that imposes an expected cost  $k > \pi_A$  on criminal agents and an expected cost  $k_0 \in (0, k]$  on non-criminal agents.

**Lemma OA.3.** *If the monitor acts as proposer when the agent is non-criminal, she demands a bribe  $\tau = k_0$  if her type is  $\eta < k_0$ , and she demands no bribe (i.e.  $\tau = 0$ ) if her type is  $\eta \geq k_0$ . A non-criminal agent accepts any offer  $\tau \leq k_0$ .*

**Proof.** Suppose the monitor makes an offer  $\tau$  to a non-criminal agent and commits to sending a false message if her proposal is rejected. In this case, it is optimal for a non-criminal agent to accept the offer if and only if  $\tau \leq k_0$ : her payoff from accepting such an offer is  $-\tau$ , while her payoff from rejecting the offer is  $-k_0$ . The monitor's payoff from making an offer

$\tau \in (0, k_0]$  is  $\tau - \eta$ , while her payoff from not demanding a bribe is 0. A type  $\eta$  monitor finds it optimal to make an offer  $\tau = k_0$  if only if  $\eta < k_0$ . ■

Lemma OA.3 implies that the payoff of a non-criminal agent is  $-(1 - \lambda)k_0F_\eta(k_0)$ . On the other hand, by the same arguments as in Section 4, the payoff of a criminal agent of type  $\pi_A$  is  $\pi_A - k + \lambda \max_\tau(k - \tau)\text{prob}(qw + \eta < \tau)$ . Therefore, when the monitor can commit to sending a false report, an agent of type  $\pi_A$  will take action  $c = 0$  if only if

$$\pi_A - (k - (1 - \lambda)k_0F_\eta(k_0)) + \lambda \max_{\tau \in [0, k]}(k - \tau)\text{prob}(qw + \eta < \tau) \leq 0.$$

From the principal's perspective, the possibility of extortion by the monitor reduces the effective punishment cost that a criminal agent incurs when the monitor sends report  $m = 1$  to  $k - (1 - \lambda)k_0F_\eta(k_0)$ . Note that this term does not depend on the distribution of wages. Hence, all the results in Sections 4 and 5 continue to hold when the monitor can commit to sending a false message.

### OA.3 Efficient contracting between the principal and monitor

Throughout the paper we assume that the principal compensates the monitor with an efficiency wage contract. This appendix shows that random incentives continue to be valuable when we allow for arbitrary contracts. We consider the same environment as in Section 3, with one minor modification: we impose a participation constraint that the agent's payoff cannot be negative. We stress, however, that the results in the main text would remain unchanged if we added this constraint.<sup>4</sup> We also assume that the cost  $k_0$  that a non-criminal agent expects from the judiciary is strictly positive.<sup>5</sup>

---

<sup>4</sup>Indeed, when the monitor is compensated with an efficiency wage  $w \geq 0$  the agent can guarantee herself a payoff of 0 by taking action  $c = 0$ . When we allow for arbitrary contracts, the agent's participation constraint rules out wage structures under which the agent needs to bribe the monitor to get a favorable report after taking action  $c = 0$ .

<sup>5</sup>As in the main text, we assume that the probability  $q$  of detecting a false report of  $m = 0$  when the agent took action  $c = 1$  is the same as the probability of detecting a false report  $m = 1$  when the agent took

Let  $s \in \{\emptyset, f\}$  denote the signal that the principal observes by scrutinizing the monitor's report: the principal observes signal  $s = f$  when she detects that the monitor's report is false, and observes signal  $s = \emptyset$  otherwise.<sup>6</sup> The principal offers a wage contract  $w(m, s)$  to the monitor, which determines the monitor's compensation as a function of the report she sends and the principal's signal. By limited liability,  $w(m, s) \geq 0$  for all  $(m, s) \in \{0, 1\} \times \{\emptyset, f\}$ .

We begin by noting that a monitor who is compensated with contract  $w(m, s)$  accepts a bribe  $\tau$  from a criminal agent if and only if  $\tau > w(1, \emptyset) - (1 - q)w(0, \emptyset) - qw(0, f)$ .

**Lemma OA.4.** *Let  $w(m, s)$  be a contract that induces the monitor to send message  $m = 0$  when the agent takes action  $c = 0$  and offers bribe  $\tau = 0$ . Then, it must be that  $w(0, \emptyset) \geq (1 - q)w(1, \emptyset) + qw(1, f)$ .*

**Proof.** When the agent takes action  $c = 0$  and offers bribe  $\tau = 0$ , the monitor's payoff from sending message  $m = 0$  is  $w(0, \emptyset)$ , while her payoff from sending message  $m = 1$  is  $(1 - q)w(1, \emptyset) + qw(1, f)$ . The monitor sends message  $m = c = 0$  if and only if  $w(0, \emptyset) \geq (1 - q)w(1, \emptyset) + qw(1, f)$ . ■

**Lemma OA.5.** *Under an optimal incentive scheme (either deterministic or random), a principal who wants to induce the agent to take action  $c = 0$  offers the monitor contracts  $w(m, s)$  with  $w(0, \emptyset) = (1 - q)w(1, \emptyset)$  and  $w(m, f) = 0$  for  $m = 0, 1$ .*

**Proof.** Suppose the incentive scheme induces the agent to take action  $c = 0$  and satisfies the agent's participation constraint. By Lemma OA.4, any contract  $w(m, s)$  that the principal offers to the monitor with positive probability must satisfy  $w(0, \emptyset) \geq (1 - q)w(1, \emptyset) + qw(1, f)$ ; otherwise the agent's expected payoff from action  $c = 0$  would be strictly negative, either because with positive probability the monitor sends a false report  $m = 1$ , or because the

---

action  $c = 0$ . Our results remain qualitatively unchanged if we allow these two probabilities to be different.

<sup>6</sup>When the monitor sends report  $m \neq c$ , the principal observes signal  $s = f$  with probability  $q$  and signal  $s = \emptyset$  with probability  $1 - q$ . When the monitor sends report  $m = c$ , the principal observes signal  $s = \emptyset$  with probability 1.



agent needs to bribe the monitor for a report  $m = 0$ . In either case, this would violate the agent's participation constraint.

This implies that under an optimal incentive scheme that induces the agent to take action  $c = 0$ , on the equilibrium path the monitor sends report  $m = 0$  and receives a wage  $w(0, \emptyset)$ . If  $w(0, \emptyset) > (1 - q)w(1, \emptyset) + qw(1, f)$  for some contract  $w(m, s)$  that is offered with positive probability, the principal would be strictly better-off by reducing  $w(0, \emptyset)$  as this would reduce wage payments and would also increase the cost of bribing the monitor.

By limited liability it must be that  $w(m, f) \geq 0$  for  $m = 0, 1$ . Setting  $w(0, f) = 0$  is optimal as it increases the cost of bribing the monitor. Finally, since  $w(0, \emptyset) = (1 - q)w(1, \emptyset) + qw(1, f)$ , setting  $w(1, f) = 0$  reduces the wage  $w(0, \emptyset)$  that the principal pays on the equilibrium path and also increases the cost of bribing the monitor. ■

We now consider the case in which the principal compensates the agent with a deterministic contract  $w(m, s)$ . The following result generalizes Lemma 1 to the current setting.

**Lemma OA.6.** *Suppose the principal uses a deterministic contract  $w(m, s)$ . Under collusion, the minimum cost of wages needed to induce the agent to be non-criminal is equal to  $\frac{1-q}{2-q} \frac{\pi_A}{q}$ .*

**Proof.** A monitor with contract  $w(m, s)$  accepts a bribe  $\tau$  from a criminal agent if and only if  $\tau > w(1, \emptyset) - (1 - q)w(0, \emptyset) - qw(0, f) = w(1, \emptyset) - (1 - q)w(0, \emptyset)$ , where the equality follows from OA.5. The agent's payoff from taking action  $c = 1$  is then  $\pi_A - \min\{k, w(1, \emptyset) - (1 - q)w(0, \emptyset)\}$ , while her payoff from taking action  $c = 0$  is 0. To induce the agent to take action  $c = 0$ , it must be that  $w(1, \emptyset) - (1 - q)w(0, \emptyset) \geq \pi_A$ . By Lemma OA.5,  $w(0, \emptyset) = (1 - q)w(1, \emptyset)$ , so the previous inequality yields  $w(0, \emptyset) \geq \frac{1-q}{2-q} \frac{\pi_A}{q}$ . ■

Consider next the case in which the principal randomizes over the monitor's contract  $w(m, s)$ . By Lemma OA.5, it is optimal for the principal to offer contracts  $w(m, s)$  such that  $w(0, \emptyset) = (1 - q)w(1, \emptyset)$  and  $w(m, f) = 0$  for  $m = 0, 1$ . Therefore, it is without loss

of optimality to focus on distributions over wages  $w(0, \emptyset)$ , with the understanding that a contract with  $w(0, \emptyset) = w \geq 0$  has  $w(1, \emptyset) = \frac{w}{1-q}$  and  $w(m, f) = 0$  for  $m = 0, 1$ .

The following result generalizes Proposition 1 to the current setting.

**Proposition OA.4.** *Under collusion, it is optimal for the principal to use random contracts. The cost-minimizing distribution  $\hat{F}_w^*$  over wages  $w(0, \emptyset)$  that induces the agent to be non-criminal is described by*

$$\forall w \in \left[ 0, \frac{\pi_A}{q} \frac{1-q}{2-q} \right], \quad \hat{F}_w^*(w) = \frac{k - \pi_A}{k - qw \frac{2-q}{1-q}}. \quad (\text{O8})$$

The corresponding cost of wages  $\hat{W}^*(\pi_A) \equiv \mathbb{E}_{\hat{F}^*}[w]$  is

$$\hat{W}^*(\pi_A) = \frac{1-q}{2-q} \frac{\pi_A}{q} \left[ 1 - \frac{k - \pi_A}{\pi_A} \log \left( 1 + \frac{\pi_A}{k - \pi_A} \right) \right] < \frac{1-q}{2-q} \frac{\pi_A}{q} \frac{\pi_A}{k}. \quad (\text{O9})$$

**Proof.** By our arguments above, a monitor with contract  $w(m, s)$  accepts a bribe  $\tau$  from a criminal agent if and only if  $\tau > w(1, \emptyset) - (1-q)w(0, \emptyset) - qw(0, f) = \frac{2-q}{1-q}qw(0, \emptyset)$ , where the last equality follows since  $w(1, \emptyset) = \frac{w(0, \emptyset)}{1-q}$  and  $w(m, f) = 0$  for  $m = 0, 1$  (Lemma OA.5). A distribution  $F$  over wages  $w(0, \emptyset)$  induces the agent to take action  $c = 0$  if and only if, for every bribe offer  $\tau \geq 0$ ,  $\pi_A - k + (k - \tau)\text{prob}(\tau > \frac{2-q}{1-q}qw) \leq 0$ , or equivalently, if and only if, for every  $\tau \geq 0$ ,  $F\left(\frac{\tau}{q} \frac{1-q}{2-q}\right) \leq \frac{k - \pi_A}{k - \tau}$ . Using the change in variable  $w = \frac{\tau}{q} \frac{1-q}{2-q}$ , we obtain that wage distribution  $F$  induces the agent to take action  $c = 0$  if and only if,

$$\forall w \in \left[ 0, \frac{\pi_A}{q} \frac{1-q}{2-q} \right], \quad F(w) \leq \frac{k - \pi_A}{k - qw \frac{2-q}{1-q}}. \quad (\text{O10})$$

By first-order stochastic dominance, it follows that in order to minimize expected wages, the optimal distribution must satisfy (O10) with equality. This implies that the optimal wage distribution is described by (O8). Expected cost expression (O9) follows from integration and straightforward computations. ■

## OA.4 Arbitrary bargaining

The model of Sections 3 and 4 simplifies the side-contracting stage by assuming take-it-or-leave-it offers. This appendix shows that random wages remain valuable under arbitrary bargaining mechanisms. We study a model in which the monitor and the agent can use any individually rational and incentive compatible mechanism at the side-contracting stage, but that is otherwise identical to the basic model in Section 3.

By the revelation principle, we can restrict attention to mechanisms under which the monitor announces her private information (i.e. her wage) and this announcement determines the bargaining outcome. Such a bargaining mechanism is characterized by two functions: (i)  $P(w)$ , the probability with which monitor and agent reach an agreement when the monitor's wage is  $w$ ; and (ii)  $\tau(w)$ , the expected transfer from the agent to the monitor when the monitor's wage is  $w$ . The monitor commits to send message  $m = 0$  if there is an agreement. If there is no agreement, the monitor sends the message that maximizes her final payoff (i.e., she sends a truthful message).

Given a wage schedule  $F$  and a mechanism  $(P, \tau)$ , the agent's expected payoff from crime is  $U_A = \pi_A - k + \int (P(w)k - \tau(w)) dF(w)$ . The individual rationality constraint of a criminal agent is  $U_A \geq \pi_A - k$ , since a criminal agent can guarantee  $\pi_A - k$  by not participating in the mechanism.

The payoff that a monitor with wage  $w$  who announces wage  $w'$  gets under mechanism  $(P, \tau)$  when the agent is criminal is  $\tilde{U}_M(w, w') = \tau(w') + (1 - P(w')q)w$ . By incentive compatibility,  $U_M(w) \equiv \tilde{U}_M(w, w) \geq \tilde{U}_M(w, w')$  for all  $w' \neq w$ . By individual rationality,  $U_M(w) \geq w$  for all  $w$ , since a monitor with wage  $w$  obtains a payoff of  $w$  by not participating in the mechanism and sending a truthful report.

Given a mechanism  $(P, \tau)$  and a wage distribution  $F$ , the weighted sum of the agent's and monitor's payoff when the agent is criminal is

$$(1 - \lambda) \int U_M(w) dF(w) + \lambda U_A, \tag{O11}$$

where the weight  $\lambda \in [0, 1]$  represents the monitor's bargaining power. For every wage schedule  $F$  and every  $\lambda \in [0, 1]$ , let  $\Gamma(F, \lambda)$  be the set of incentive compatible and individually rational bargaining mechanisms that maximize (O11). We assume that, at the side-contracting stage, the monitor and the agent use a bargaining mechanism in  $\Gamma(F, \lambda)$ . Let  $\tilde{U}_A(F, \lambda)$  be the lowest utility that a criminal agent gets under a bargaining mechanism in  $\Gamma(F, \lambda)$ . The agent has an incentive to be non-criminal if  $\tilde{U}_A(F, \lambda) \leq 0$ .

The following result generalizes Proposition 1 to this setting.

**Proposition OA.5.** *Suppose that, at the collusion stage, the monitor and the agent use an incentive compatible and individually rational mechanism that maximizes (O11).*

(i) *If  $\lambda \in (1/2, 1]$ , the cost minimizing wage distribution  $\tilde{F}_w^*$  that induces the agent to be non-criminal is described by*

$$\forall w \in [0, \pi_A/q], \quad \tilde{F}_w^*(w) = \left( \frac{k - \pi_A}{k - qw} \right)^{\frac{2\lambda-1}{\lambda}}. \quad (\text{O12})$$

(ii) *If  $\lambda \in [0, 1/2]$ , the cost minimizing wage distribution  $\tilde{F}_w^*$  that induces the agent to be non-criminal has  $\tilde{F}_w^*(0) = 1$ .*

**Proof.** By standard arguments, any incentive compatible mechanism  $(P, \tau)$  must satisfy: (i)  $P(w)$  is decreasing, and (ii)  $U'_M(w) = 1 - qP(w)$  a.e.. This last condition and the monitor's individual rationality constraint (i.e.,  $U_M(w) \geq w$  for all  $w$ ) imply that  $U_M(w) = \int_w^{\bar{w}} qP(\tilde{w})d\tilde{w} + w + c$  for some constant  $c \geq 0$  (where  $\bar{w}$  is the highest wage in the support of  $F$ ). Since  $U_M(w) = \tau(w) + (1 - qP(w))w$ ,  $\tau(w) = P(w)qw + \int_w^{\bar{w}} qP(\tilde{w})d\tilde{w} + c$ . The weighted

sum of payoffs when the agent is criminal is

$$\begin{aligned}
& (1 - \lambda) \int_{\underline{w}}^{\bar{w}} U_M(w) dF(w) + \lambda U_A \\
&= \int_{\underline{w}}^{\bar{w}} [(1 - \lambda)(\tau(w) + (1 - qP(w))w) + \lambda(P(w)k - \tau(w))] dF(w) + \lambda(\pi_A - k) \\
&= \int_{\underline{w}}^{\bar{w}} [P(w)\lambda(k - qw) + (1 - \lambda)w] dF(w) + \lambda(\pi_A - k) + (1 - 2\lambda) \left( \int_{\underline{w}}^{\bar{w}} qP(w)F(w)dw + c \right).
\end{aligned} \tag{O13}$$

We use the following lemma.

**Lemma OA.7.** *For all  $\lambda \in (1/2, 1]$ , the mechanism  $(P, \tau)$  that maximizes (O13) has: (i)  $P(w) = 1$  if  $w < w^*$  and  $P(w) = 0$  if  $w > w^*$  for some  $w^* \in [\underline{w}, \bar{w}]$ , and (ii)  $\tau(w) = P(w)qw + \int_w^{\bar{w}} qP(\tilde{w})d\tilde{w}$ .*

**Proof.** Note first that (O13) is maximized by setting  $c = 0$  when  $\lambda \in (1/2, 1]$ . Moreover, when  $\lambda \in (1/2, 1]$  any mechanism  $(P, \tau)$  that maximizes (O13) must be such  $P(w) = 0$  for all  $w \geq k/q$ .

We now show that the mechanism that maximize (O13) is such that  $P(w)$  only takes values 0 or 1. From above, we know that  $P(w) = 0$  for all  $w \geq k/q$ . Suppose by contradiction that there exists an interval  $V \subset [0, k/q]$  such that  $P(w) \in (0, 1)$  for all  $w \in V$ , and let  $H \equiv \int_V \lambda(k - qw)dF(w) + (1 - 2\lambda) \int_V qF(w)dw$ . If  $H \geq 0$ , increasing  $P(w)$  over this interval (subject to the constraint that  $P$  is decreasing) makes (O13) larger. If  $H < 0$ , decreasing  $P(w)$  over this interval (subject to the constraint that  $P$  is decreasing) also makes (O13) larger. Such improvements are exhausted when  $P(w)$  only takes values 0 and 1.<sup>7</sup> Since  $P(\cdot)$  is decreasing, when  $P(\cdot)$  only takes values 0 or 1 there must exist a wage  $w^*$  such that  $P(w) = 1$  if  $w < w^*$  and  $P(w) = 0$  if  $w > w^*$ . Finally, since (O13) is maximized by setting

---

<sup>7</sup>Note that these changes in  $P(w)$  do not conflict with the participation constraints of monitor and agent. Indeed,  $U_M(w) = \int_w^{\bar{w}} qP(\tilde{w})d\tilde{w} + w \geq w$  for any incentive compatible mechanism  $(P, \tau)$ . Moreover, for all  $w$ ,  $\tau(w) = P(w)qw + \int_w^{\bar{w}} qP(\tilde{w})d\tilde{w} \leq P(w)k$ , where the inequality follows since any mechanism that maximizes (O13) has  $P(w) = 0$  for all  $w \geq k/q$  and since  $P(\cdot)$  is decreasing. Hence,  $U_A = \pi_A - k + \int (P(w)k - \tau(w))dF(w) \geq \pi_A - k$ .

$c = 0$  when  $\lambda \in (1/2, 1]$ ,  $\tau(w) = P(w)qw + \int_w^{\bar{w}} qP(\tilde{w})d\tilde{w}$ . Since  $P(w) = 1$  if  $w < w^*$  and  $P(w) = 0$  if  $w > w^*$ , it follows that  $\tau(w) = qw^*$  if  $w < w^*$  and  $\tau(w) = 0$  if  $w > w^*$ . ■

We now conclude the proof of Proposition OA.5, beginning with point (i). Fix  $\lambda \in (1/2, 1]$  and let  $F$  be a cost-minimizing wage schedule that induces the agent to be non-criminal. Let  $(P, \tau)$  be the mechanism that maximizes the weighted sum of payoffs (O13) under distribution  $F$ . By Lemma OA.7,  $P(w) = \mathbf{1}_{w \leq w^*}$  and  $\tau(w) = qw^*\mathbf{1}_{w \leq w^*}$  for some  $w^*$ . Under this mechanism (O13) becomes

$$\begin{aligned} & \lambda \left[ F(w^*)k - \int_0^{w^*} qwdF(w) + \pi_A - k \right] + (1 - \lambda) \int w dF(w) + (1 - 2\lambda) \int_0^{w^*} qF(w)dw \\ &= \lambda [F(w^*)(k - qw^*) + \pi_A - k] + (1 - \lambda) \int w dF(w) + (1 - \lambda) \int_0^{w^*} qF(w)dw, \end{aligned}$$

where we used  $\int_0^{w^*} qwdF(w) = qw^*F(w^*) - \int_0^{w^*} qF(w)dw$ . Since  $(P, \tau)$  maximizes the weighted sum of payoffs, for all  $\hat{w} \neq w^*$  it must be that

$$\lambda F(w^*)(k - qw^*) + (1 - \lambda) \int_0^{w^*} qF(w)dw \geq \lambda F(\hat{w})(k - q\hat{w}) + (1 - \lambda) \int_0^{\hat{w}} qF(w)dw$$

Otherwise, if the inequality did not hold for some  $\hat{w} \neq w^*$ , the weighted sum of payoffs would be strictly larger under mechanism  $(\hat{P}, \hat{\tau})$  with  $\hat{P}(w) = 1$  if  $w < \hat{w}$  and  $\hat{P}(w) = 0$  if  $w > \hat{w}$ .

For any  $\hat{w} \in \text{supp } F$ , let  $(P_{\hat{w}}, \tau_{\hat{w}})$  be the mechanism with  $P_{\hat{w}}(w) = \mathbf{1}_{\{w \leq \hat{w}\}}$  and  $\tau_{\hat{w}}(w) = \mathbf{1}_{\{w \leq \hat{w}\}}q\hat{w}$ . Recall that  $\Gamma(F, \lambda)$  is the set of bargaining mechanisms that maximize (O13) and that  $\tilde{U}_A(F, \lambda)$  is the lowest utility that a criminal agent gets under a mechanism in  $\Gamma(F, \lambda)$ . By our arguments above,

$$\Gamma(F, \lambda) = \left\{ (P_{\hat{w}}, \tau_{\hat{w}}) : \hat{w} \in \arg \max_{w'} \lambda F(w')(k - qw') + (1 - \lambda) \int_0^{w'} qF(w)dw \right\}.$$

Suppose that there exists  $w_1$  and  $w_2 > w_1$  such that  $(P_w, \tau_w) \in \Gamma(F, \lambda)$  for  $w = w_1, w_2$ . Note that the agent's payoff from being criminal under mechanism  $(P_w, \tau_w)$  is  $F(w)(k - qw) + \pi_A - k$ .

Since  $(P_w, \tau_w) \in \Gamma(F, \lambda)$  for  $w = w_1, w_2$ ,

$$\lambda F(w_1)(k - qw_1) + (1 - \lambda) \int_0^{w_1} qF(w)dw = \lambda F(w_2)(k - qw_2) + (1 - \lambda) \int_0^{w_2} qF(w)dw$$

and so  $F(w_2)(k - qw_2) < F(w_1)(k - qw_1)$ . This implies that,  $\tilde{U}_A(F, \lambda) = F(\tilde{w})(k - q\tilde{w}) + \pi_A - k$ , where  $\tilde{w} \equiv \sup\{\hat{w} \in \text{supp } F : \hat{w} \in \arg \max_{w'} \lambda F(w')(k - qw') + (1 - \lambda) \int_0^{w'} qF(w)dw\}$ . Since  $F$  induces the agent to be non-criminal,  $\tilde{U}_A(F, \lambda) = F(\tilde{w})(k - q\tilde{w}) + \pi_A - k \leq 0$ .

Let  $\bar{w}$  be the highest wage in the support of  $F$ . We now show that, if  $F$  is an optimal distribution, it must be that  $\bar{w} \in \arg \max_{w'} \lambda F(w')(k - qw') + (1 - \lambda) \int_0^{w'} qF(w)dw$ . Suppose by contradiction that this is not true, so that  $\bar{w} > \tilde{w} = \sup\{\hat{w} \in \text{supp } F : \hat{w} \in \arg \max_{w'} \lambda F(w')(k - qw') + (1 - \lambda) \int_0^{w'} qF(w)dw\}$ . Pick  $\epsilon \in (0, \bar{w} - \tilde{w})$  small and let  $F^\epsilon$  be a c.d.f. with  $F^\epsilon(w) = F(w)$  for all  $w < \bar{w} - \epsilon$  and  $F^\epsilon(\bar{w} - \epsilon) = 1$ . By first-order stochastic dominance,  $\mathbb{E}_{F^\epsilon}[w] < \mathbb{E}_F[w]$ . By the definition of  $\tilde{w}$ ,

$$\lambda F(\tilde{w})(k - q\tilde{w}) + (1 - \lambda) \int_0^{\tilde{w}} qF(w)dw \geq \lambda F(\hat{w})(k - q\hat{w}) + (1 - \lambda) \int_0^{\hat{w}} qF(w)dw,$$

for all  $\hat{w}$ , with strict inequality for all  $\hat{w} \in (\tilde{w}, \bar{w}]$ . Therefore, there exists  $\epsilon > 0$  small enough such that, for all  $\hat{w}$ ,

$$\lambda F^\epsilon(\tilde{w})(k - q\tilde{w}) + (1 - \lambda) \int_0^{\tilde{w}} qF^\epsilon(w)dw \geq \lambda F^\epsilon(\hat{w})(k - q\hat{w}) + (1 - \lambda) \int_0^{\hat{w}} qF^\epsilon(w)dw$$

This implies that mechanism  $(P_{\tilde{w}}, \tau_{\tilde{w}})$  is still optimal under distribution  $F^\epsilon$ , and so  $\tilde{U}_A(F^\epsilon, \lambda) \leq F(\tilde{w})(k - q\tilde{w}) + \pi_A - k \leq 0$ . But this cannot be, since  $F$  is a cost-minimizing distribution that induces the agent to be non-criminal. Therefore, if  $F$  is optimal it must be that  $\bar{w} = \sup\{\hat{w} \in \text{supp } F : \hat{w} \in \arg \max_{w'} \lambda F(w')(k - qw') + (1 - \lambda) \int_0^{w'} qF(w)dw\}$ . The agent's payoff from being criminal under mechanism  $(P_{\bar{w}}, \tau_{\bar{w}})$  is  $k - q\bar{w} + \pi_A - k \leq 0 \iff \bar{w} \geq \frac{\pi_A}{q}$ .

By the arguments above, for all  $\hat{w} \in [0, \bar{w}]$ ,

$$\begin{aligned} \lambda(k - q\bar{w}) + (1 - \lambda) \int_0^{\bar{w}} qF(w)dw &\geq \lambda F(\hat{w})(k - q\hat{w}) + (1 - \lambda) \int_0^{\hat{w}} qF(w)dw \\ \iff \lambda(k - q\bar{w}) + (1 - \lambda) \int_{\hat{w}}^{\bar{w}} qF(w)dw &\geq \lambda F(\hat{w})(k - q\hat{w}) \end{aligned} \quad (\text{O14})$$

We now show that, if  $F$  is an optimal distribution, (O14) must hold with equality for all  $\hat{w} \in [0, \bar{w}]$ . Suppose by contradiction that there is an interval  $[w_1, w_2] \subset [0, \bar{w})$  such that (O14) is slack for all  $\hat{w} \in [w_1, w_2]$ . By first-order stochastic dominance, increasing  $F(\cdot)$  over  $[w_1, w_2]$  (subject to the constraint that  $F$  is increasing) reduces expected wage payments. Moreover, increasing  $F(\cdot)$  over  $[w_1, w_2]$  relaxes (O14) for all  $\hat{w} < w_1$  and does not affect (O14) for all  $\hat{w} > w_2$ . This implies that mechanism  $(P_{\bar{w}}, \tau_{\bar{w}})$  still maximizes the weighted sum of payoffs (O13) after increasing  $F(\cdot)$  slightly over  $[w_1, w_2]$ , and so the agent's payoff from being criminal is  $k - q\bar{w} + \pi_A - k \leq 0$ . But this cannot be, since  $F$  is a cost-minimizing distribution that induces the agent to be non-criminal. Therefore, if  $F$  is optimal, (O14) must hold with equality for all  $\hat{w} \leq \bar{w}$ .

Since (O14) holds with equality for all  $\hat{w} \leq \bar{w}$ ,  $\lambda F(\hat{w})(k - q\hat{w}) + (1 - \lambda) \int_0^{\hat{w}} qF(w)dw$  is constant over  $[0, \bar{w}]$ . Differentiating this expression with respect to  $\hat{w}$ , it must be that

$$F'(\hat{w})\lambda[k - q\hat{w}] + qF(\hat{w})(1 - 2\lambda) = 0. \quad (\text{O15})$$

The solution to the differential equation (O15) is  $F(w) = C \left( \frac{1}{k - qw} \right)^{\frac{2\lambda - 1}{\lambda}}$ , where  $C$  is a constant such that  $F(\bar{w}) = 1$ ; i.e.,  $C = (k - q\bar{w})^{\frac{2\lambda - 1}{\lambda}}$ . Finally, by our arguments above, under distribution  $F$  the agent will have an incentive to be non-criminal as long as  $k - q\bar{w} + \pi_A - k \leq 0 \iff \bar{w} \geq \frac{\pi_A}{q}$ . Since the constant  $C$  is decreasing in  $\bar{w}$ , an optimal distribution must have  $\bar{w} = \frac{\pi_A}{q}$ . Hence,  $C = (k - \pi_A)^{\frac{2\lambda - 1}{\lambda}}$ , so the optimal distribution is (O12).

We now turn to point (ii). When  $\lambda \leq 1/2$ , the mechanism  $(P, \tau)$  that maximizes (O13) must make the constant  $c$  as large as possible, subject to the agent's IR constraint; that is, subject to  $\pi_A - k + \int [P(w)k - \tau(w)]dF(w) \geq \pi_A - k$ . Recall that  $\tau(w) = P(w)qw +$



$\int_w^{\bar{w}} qP(\tilde{w})d\tilde{w}+c$ . The maximum is achieved by choosing  $c$  such that  $\int [P(w)k-\tau(w)]dF(w) = 0$ . Therefore, for  $\lambda \leq 1/2$  the agent's payoff from engaging in crime under a mechanism that maximizes (O13) is  $\pi_A - k < 0$ , regardless of the wage schedule. This implies that the agent has an incentive to be non-criminal even when  $F$  has all its mass at  $w = 0$ . ■

We end this appendix by noting that the results above generalize to settings in which the agent is privately informed about the benefit  $\pi_A$  from crime. Given a wage profile  $F_w$ , the payoff an agent of type  $\pi_A$  gets from taking action  $c = 1$  is  $U_A(\pi_A) = \pi_A - k + \int (P(w; F_w)k - \tau(w; F_w)) dF(w)$ , where  $(P(w), \tau(w))$  is the mechanism that maximizes the weighted sum of payoffs (O13).<sup>8</sup> Since  $U_A(\pi_A)$  is increasing in  $\pi_A$ , agents follow a threshold strategy: for any wage schedule  $F_w$ , there is a cutoff  $\bar{\pi}_A(F_w)$  such that an agent of type  $\pi_A$  is criminal if and only if  $\pi_A > \bar{\pi}_A(F_w)$ . For any cutoff  $\pi_A$ , Proposition OA.5 characterizes the cheapest wage distribution that attains this cutoff.

---

<sup>8</sup>Note that, given wage profile  $F_w$ , the mechanism  $(P(w), \tau(w))$  that maximizes (O13) is independent of the agent's type  $\pi_A$ .