

# But Who will Monitor the Monitor?\*

David Rahman<sup>†</sup>  
University of Minnesota

January 16, 2009.

## Abstract

Consider a group of individuals in a strategic environment with moral hazard and adverse selection, and suppose that providing incentives for a given outcome requires a monitor to detect deviations. What about the monitors' deviations? This paper analyzes mediated contracts to answer such a question, and asserts that the monitor's deviations are effectively irrelevant. Hence, nobody needs to monitor the monitor. The paper also characterizes exactly when such contracts can provide monitors with the right incentives to perform. In doing so, several characterizations of virtual implementation are derived.

*JEL Classification:* D21, D23, D82.

*Keywords:* contracts, private monitoring, communication equilibrium.

---

\*This paper subsumes [Rahman \(2005a\)](#), which in turn is based on Chapter 3 of [Rahman \(2005b\)](#).

<sup>†</sup>Financial support from the Spanish Ministry of Education's Research Grant No. SEJ 2004-07861 while at Universidad Carlos III de Madrid is gratefully acknowledged.

# 1 Introduction

Ann owns a restaurant. She hires Bob to tally the till every night and report back any mismatch between the till and that night’s bills. Ann can motivate Bob to exert such effort and report truthfully any mismatch by sometimes secretly taking money from the till herself and offering him the following incentive scheme: if Ann took some money, she will pay Bob his wage only when he reports a mismatch; if Ann did not take any money, she will pay Bob only when a mismatch is not reported.

Bob faces a secret contract: his report-contingent wage is unknown to him a priori (it depends on whether or not Ann secretly took some money). If Bob fails to exert effort, he won’t know what to report in order to secure his wage. However, if he does his job he’ll discover whether or not there is a mismatch and deduce from this Ann’s behavior. Only after tallying the till will Bob know what to report in order to receive his wage, which turns out to be optimally truthful.

This paper studies contracts like Bob’s<sup>1</sup> and how they might help organizations to function productively. By allocating different information to team members, secret contracts often provide better incentives to perform with an intuitive organizational design. Thus, they give Bob incentives to acquire costly information and reveal it, and provide Ann with enough a priori knowledge to distinguish working from shirking. In general, they provide a way to “monitor the monitor” ([Section 2.1](#)).

Consider a hypothetical organization whose individuals are subject to moral hazard but with rich communication protocols: access to (i) a disinterested mediator or machine that makes confidential, verifiable but non-binding recommendations to agents, and (ii) (linear) “money” transfers that may depend on the mediator’s recommendations and individual reports (such as Bob’s). A contract thus involves instructions and payments:, i.e., a way of telling people what to do and a way of rewarding them.

When can this organization actually overcome moral hazard with secret contracts? Below, we study incentives in such a team to answer this question in various contexts.

---

<sup>1</sup>These contracts are pervasive. For instance, TSA screeners are evaluated with “covert testing” ([TSA, 2004](#), p. 5); police use young ‘drinkers’ to ensure that bartenders check IDs ([Cheslow, 2005](#)).

Formally, we consider contractual arrangements subject to incentive compatibility as described by Myerson’s (1986) *communication equilibrium*. We obtain minimal conditions on a team’s primitives—its “monitoring technology” (what Bob can see when) and individual preferences (whether Bob prefers to work or shirk)—such that incentive compatibility is not a binding constraint for the team.

[Theorem 1](#) provides a necessary and sufficient condition on a monitoring technology—called detecting unilateral disobedience (DUD)—for every team outcome (e.g., Bob works) to be approximately enforceable, i.e., an incentive compatible outcome exists arbitrarily close to it. DUD requires that every disobedient deviation by any individual be statistically detectable with some reaction by others, although different deviations may be detected with different reactions. This key property distinguishes DUD substantively from the literature<sup>2</sup> ([Section 3.1](#) has a detailed literature review). Therefore, DUD is a weak restriction. It is also generic ([Theorem 2](#)).

Secret contracts add value *not* by approximate enforcement ([Corollaries 1 and 2](#)),<sup>3</sup> but by allowing “monitors” to *follow* “deviators” in a hypothetical game of hide and seek, even though in fact they move simultaneously. To illustrate, suppose Bob shirks. If he also reports no mismatch then Ann can hypothetically “react” by secretly taking some money to prove him wrong, whereas if he reports a mismatch then Ann can choose not to take any money. By [Theorem 1](#), such disobedience (e.g., Bob shirking) is detectable in this sense if and only if obedience is enforceable with secret contracts.

[Theorems 3 and 5](#) ([Section 4](#)) extend these results in two important directions that help to clarify the differences between exact and approximate enforcement.

[Theorem 3](#) characterizes monitoring technologies that approximately enforce a fixed outcome rather than every outcome simultaneously ([Theorem 1](#)), regardless of individual preferences. Interestingly, [Theorem 3](#) reconciles an infinite regress inherent in monitoring. Suppose that providing incentives for a given outcome requires a monitor to detect deviations. What about the monitor’s deviations? [Theorem 3](#) answers this question by asserting that effectively the monitor’s deviations are irrelevant. Indeed, if they are detectable then they can be easily discouraged with contingent payments. Otherwise, if the monitor’s deviations are undetectable then

---

<sup>2</sup>For instance, conditions like *individual full rank* of [Fudenberg et al. \(1994\)](#) require that every deviation be detected by the same “reaction,” making it more difficult to detect deviations.

<sup>3</sup>Even though we sometimes rely on approximation to expand contractual possibilities, our key insight is the use of *mediated transfers*—not approximating outcomes—to provide incentives.

the deviations themselves still detect others’ deviations from the given outcome, and so they continue to fulfill the required monitoring role. Evidently, this argument also applies to the monitor’s deviations from these deviations, and so forth. [Theorem 3](#) reconciles this infinite regress by showing that under standard conditions (e.g., in a finite game) not every behavior by the monitor can have a profitable, undetectable deviation. Therefore, to approximately enforce an arbitrary outcome with infrequent monitoring, every deviation from the outcome must in principle be detectable with some monitoring behavior, but deviations away from the monitoring behavior itself need not be detectable. Heuristically, *nobody needs to monitor the monitor*.

[Theorem 5](#) extends [Theorem 3](#) by fixing individual preferences and finding joint conditions on preferences and the monitoring technology that characterize approximate enforcement. Intuitively, profitable deviations must be discouraged “uniformly” and “credibly.” Uniform detection allows for infinitesimal deviations to be discouraged even if they are only infinitesimally detectable ([Example 5](#)). Credibility is necessary when deviations are discouraged with the use of others’ actions rather than with contingent payments, and this disciplining behavior must be incentive compatible.

The paper is organized as follows. [Section 2](#) presents two motivating examples that guide our main results. [Section 3](#) develops the model. [Section 3.1](#) defines DUD, characterizes its incentive properties, finds conditions for its generic satisfaction, and relates it to the literature. [Section 4](#) extends the model by characterizing exact and approximate enforcement of fixed outcomes with and without fixed preferences, and accommodates complications such as participation constraints, limited liability, and even coalitional deviations. [Section 5](#) concludes. Omitted proofs and ancillary results appear in [Appendices A](#) and [B](#).

## 2 Example

We begin our analysis of secret contracts with two leading examples that attempt to capture the intuition behind our main results, [Theorems 1](#) and [3](#). The first example considers an environment that typifies the strategic interaction between a principal, a worker, and a monitor. The second example suggests an intuitive way of attaining approximately efficient partnerships with budget balance.

## 2.1 Robinson and Friday

There are two agents: Robinson, who can either monitor or rest, and Friday, who can either work or shirk. A *mediating* principal makes possibly secret recommendations to the agents and enforces contingent contractual payments. Robinson (the row player) and Friday (the column player) interact according to the left bi-matrix below.

	work	shirk
monitor	2, -1	-1, 0
rest	3, -1	0, 0

Utility Payoffs

	work	shirk
monitor	1, 0	0, 1
rest	1/2, 1/2	1/2, 1/2

Signal Probabilities

There are two signals,  $g$  and  $b$ , on which to condition linear transfers. Their conditional probability is given in the right bi-matrix above. In words, if Robinson monitors he observes Friday's effort, whereas if he rests then the signal is uninformative.

Although clearly the efficient profile (rest,work) is unenforceable, we can get arbitrarily close even if only Robinson observes the signal and it is not verifiable.<sup>4</sup> For the principal to write signal-contingent contracts, he must first solicit the realizations from Robinson, who may in principle misreport them.<sup>5</sup> We approximate (rest,work) by having Friday mix between working and shirking and Robinson's report-contingent payments depend on Friday's recommendation, thereby "monitoring the monitor." Specifically, the following correlated strategy is incentive compatible given  $\mu \in (0, 1)$ :  
(i) Robinson is told to monitor with probability  $\sigma$  (and rest with probability  $1 - \sigma$ ),  
(ii) Friday is independently told to work with probability  $\mu$  (to shirk with  $1 - \mu$ ), and  
(iii) the principal enforces the following *secret contract*:

	(monitor,work)	(monitor,shirk)	(rest,work)	(rest,shirk)
$g$	$1/\mu, 1/\sigma$	0, 0	0, 0	0, 0
$b$	0, 0	$1/(1 - \mu), 0$	0, 0	0, 0

The table reads as follows. The leftmost column says that Robinson is paid  $\$1/\mu$  if he reports  $g$  and  $\$0$  if  $b$  when (monitor,work) was recommended, whereas Friday is paid

<sup>4</sup>If signals are publicly verifiable, the correlated strategy  $\sigma[(\text{monitor}, \text{work})] + (1 - \sigma)[(\text{rest}, \text{work})]$ , where  $[a]$  means Dirac measure for any action profile  $a$ , is enforced for all  $\sigma \in (0, 1]$  with [Holmström's \(1982\) group penalties](#), e.g., by paying Robinson  $\$2$  and Friday  $\$1/\sigma$  if  $g$  and both agents  $\$0$  if  $b$ .

<sup>5</sup>Now group penalties break down, since then Robinson reports  $g$  and rests, hence Friday shirks. Furthermore, if Robinson was paid independently of his report then although he would happily tell the truth, he would find no reason to monitor.

$\$1/\sigma$  if  $g$  is reported and  $\$0$  if  $b$ , etc. Honesty and obedience to the mediator is now incentive compatible. Letting  $\sigma \rightarrow 0$  and  $\mu \rightarrow 1$ , (rest,work) can now be approached.

Intuitively, Robinson is rewarded only when he reports  $g$  if Friday was asked to work and  $b$  if Friday was asked to shirk. Robinson, like Bob, faces a “trick question.”

Secret contracts add value in this example because they allow different correlated strategies to detect different deviation plans, unlike just signal-contingent contracts. In other words, this is as if a correlated strategy is chosen *after* agents choose deviation plans in order to detect them. To illustrate, suppose that Robinson is asked to monitor but instead chooses to rest and report  $g$ . The mediator can “react” by asking Friday to shirk, which would lead to  $b$  if Robinson monitored and reported truthfully. Similarly, if Robinson plans to rest and report  $b$  then Friday can be asked to work instead, and Robinson’s deviation is detected again.

The key idea behind [Theorem 1](#) shows that, therefore, Robinson can be dissuaded from resting. However, with only signal-contingent contracts ([Corollary 1](#)), detecting Robinson’s deviations requires the principal to fix Friday’s behavior in advance. If Friday works with fixed probability  $\mu$  then Robinson can rest and report  $g$  with probability  $\mu$ . Now Robinson can deviate without being detected, and no contract contingent only on signals can induce him to monitor.

### 3 Model

This section develops the main model of secret contracts, whose purpose is to characterize a team’s enforceable outcomes. Firstly, basic notation is introduced, the timing of interaction amongst team members is described explicitly, and several notions of enforcement are formally defined that will be used extensively later.

Section 3.1 then extrapolates from the leading example in Section 2.1. A notion of detection of deviation plans is introduced and the equivalence between detection and enforcement is derived in terms of a hypothetical zero-sum game of hide and seek where the hider (a deviator) moves first and the seeker (a monitor) moves second. Allowing the seeker to move second is shown to characterize the value of secret contracts. Formally, an outcome is shown to be enforceable if and only if the seeker wins in this hypothetical hide-and-seek game. A notion of “almost perfect monitoring” is also proposed, called detecting unilateral disobedience, and conditions are provided for it to obtain generically. Section 3.1 ends with a literature review.

We begin by defining the basic strategic environment. Let  $I = \{1, \dots, n\}$  be a finite set of agents,  $A_i$  a finite set of actions available to any agent  $i \in I$ , and  $A = \prod_i A_i$  the (nonempty) space of action profiles. Let  $v_i(a)$  denote the utility to agent  $i \in I$  from action profile  $a \in A$ . A *correlated strategy* is a probability measure  $\sigma \in \Delta(A)$ .<sup>6</sup>

Let  $S_i$  be a finite set of *private signals* observable only by individual member  $i \in I$  and  $S_0$  a finite set of *publicly verifiable* signals. Let

$$S := \prod_{j=0}^n S_j$$

be the (nonempty) product space of all observable signals. A *monitoring technology* is a measure-valued map  $\Pr : A \rightarrow \Delta(S)$ , where  $\Pr(s|a)$  stands for the conditional probability that  $s = (s_0, s_1, \dots, s_n) \in S$  was observed given that the team played  $a = (a_1, \dots, a_n) \in A$ .

Assume that the team has access to *linear transfers*. An *incentive scheme* is any map  $\zeta : I \times A \times S \rightarrow \mathbb{R}$  that assigns monetary transfers contingent on individuals, recommended actions, and *reported* signals. It is assumed that recommendations are verifiable.<sup>7</sup> Rather than focus on incentive schemes  $\zeta$ , we will also study *probability*

---

<sup>6</sup>If  $X$  is a finite set,  $\Delta(X) = \{\mu \in \mathbb{R}_+^X : \sum_x \mu(x) = 1\}$  is the set of probability vectors on  $X$ .

<sup>7</sup>If recommendations were not directly verifiable, then agents could be asked to announce theirs

weighted transfers,  $\xi : I \times A \times S \rightarrow \mathbb{R}$ . For any recommendation  $a \in A$  with  $\sigma(a) > 0$ , one may think of  $\xi$  as solving  $\xi_i(a, s) = \sigma(a)\zeta_i(a, s)$  for some  $\zeta$ . For any  $a \in A$  with  $\sigma(a) = 0$  and  $\xi(a) \neq 0$ , one may think of  $\xi$  as either arising from unbounded incentive schemes (i.e.,  $\zeta_i(a, s) = \pm\infty$ ) or as the limit of a sequence  $\{\sigma^m \zeta^m\}$ . This change of variables from  $\zeta$  to  $\xi$  is explained further in [Section 4.1](#).

The timing of team members' interaction runs as follows. Firstly, agents agree upon some *contract*  $(\sigma, \zeta)$  consisting of a correlated strategy  $\sigma$  and an incentive scheme  $\zeta$ . A profile of recommendations is drawn according to  $\sigma$  and made to agents confidentially and verifiably by some machine. Agents then simultaneously take some action. Taken actions are neither verifiable nor directly observable. Next, agents observe their unverifiable private signals and submit a verifiable report of their observations (given by a signal) before observing the public signal (not essential, just simplifying). Finally, recommendation- and report-contingent transfers are made according to  $\zeta$ .

If every agent obeys his recommendation and reports truthfully, the expected utility to agent  $i$  (before recommendations are actually made) from a contract  $(\sigma, \zeta)$  is

$$\sum_{a \in A} \sigma(a) v_i(a) - \sum_{(a, s)} \sigma(a) \zeta_i(a, s) \Pr(s|a).$$

Of course, Mr.  $i$  may disobey his recommendation  $a_i$  to play some other action  $b_i$  and lie about his privately observed signal. A *reporting strategy* is a map  $\rho_i : S_i \rightarrow S_i$ , where  $\rho_i(s_i)$  is the reported signal when Mr.  $i$  privately observes  $s_i$ . Let  $R_i$  be the set of all reporting strategies for agent  $i$ . The *truthful reporting strategy* is the identity map  $\tau_i : S_i \rightarrow S_i$  with  $\tau_i(s_i) = s_i$ . Thus, both  $\zeta_i(a, s_{-i}, \tau_i(s_i)) = \zeta_i(a, s)$  and  $\xi_i(a, s_{-i}, \tau_i(s_i)) = \xi_i(a, s)$ .<sup>8</sup> The space of pure *deviations* for  $i$  is therefore  $A_i \times R_i$ .

For every agent  $i$  and every deviation  $(b_i, \rho_i)$ , the conditional probability that signal profile  $s$  will be reported when everyone else is honest and plays  $a_{-i} \in A_{-i}$  equals

$$\Pr(s|a_{-i}, b_i, \rho_i) := \sum_{t_i \in \rho_i^{-1}(s_i)} \Pr(s_{-i}, t_i|a_{-i}, b_i).$$

When all other agents are honest and obedient, the utility to  $i$  from deviating to  $(b_i, \rho_i)$  conditional on being recommended to play  $a_i$  under contract  $(\sigma, \zeta)$  equals

$$\sum_{a_{-i}} \frac{\sigma(a)}{\sigma(a_i)} v_i(a_{-i}, b_i) - \sum_{(a_{-i}, s)} \frac{\sigma(a)}{\sigma(a_i)} \zeta_i(a, s) \Pr(s|a_{-i}, b_i, \rho_i),$$

---

as verifiable messages. However, this would involve some loss of generality ([Example 2](#)).

<sup>8</sup>We will often use the notation  $s = (s_{-i}, s_i)$  and  $a = (a_{-i}, a_i)$  for any  $i$ , where  $s_i \in S_i$  and  $s_{-i} \in S_{-i} = \prod_{j \neq i} S_j$ ; similarly for  $A_{-i}$ .



where  $\sigma(a_i) = \sum_{a_{-i}} \sigma(a) > 0$  is the probability that  $a_i$  was recommended.

A team's *metering problem* is to find a contract  $(\sigma, \zeta)$  that makes incentive compatible obeying recommended behavior as well as honest reporting of monitoring signals. This is captured by the following family of inequalities.

$$\forall i \in I, a_i \in A_i, (b_i, \rho_i) \in A_i \times R_i, \\ \sum_{a_{-i}} \sigma(a) (v_i(a_{-i}, b_i) - v_i(a)) \leq \sum_{(a_{-i}, s)} \sigma(a) \zeta_i(a, s) (\Pr(s|a_{-i}, b_i, \rho_i) - \Pr(s|a)). \quad (*)$$

The left-hand side reflects the *deviation gain* in terms of utility<sup>9</sup> for an agent  $i$  from playing  $b_i$  when asked to play  $a_i$ . The right-hand side reflects his *contractual loss* from deviating to  $(b_i, \rho_i)$  relative to honesty and obedience (i.e., playing  $a_i$  after being asked to do so and reporting according to  $\tau_i$ ). Such a loss originates from two sources. On the one hand, playing  $b_i$  instead of  $a_i$  may change conditional probabilities over signals. On the other, reporting according to  $\rho_i$  may affect conditional payments.

**Definition 1.** A correlated strategy  $\sigma$  is *exactly enforceable* (or simply *enforceable*) if there exists an incentive scheme  $\zeta : I \times A \times S \rightarrow \mathbb{R}$  to satisfy  $(*)$  for all  $(i, a_i, b_i, \rho_i)$ . A correlated strategy  $\sigma$  is *approximately enforceable* if a sequence  $\{\sigma^m\}$  of enforceable correlated strategies exists with  $\sigma^m \rightarrow \sigma$ .

A correlated strategy is approximately enforceable if it is the limit of exactly enforceable ones. E.g., in [Section 2.1](#) the correlated strategy [(rest,work)] is approximately enforceable but not enforceable.

---

<sup>9</sup>Specifically, in terms of probability weighted utility, weighted by  $\sigma(a_i)$ . If  $a_i$  is never recommended then  $\sigma(a_i) = 0$  and both sides of the inequality equal zero.

### 3.1 Detection

We now provide a notion of detection that is shown to be equivalent to enforcement. A *deviation plan* for any agent  $i$  is a map  $\alpha_i : A_i \rightarrow \Delta(A_i \times R_i)$ , where  $\alpha_i(b_i, \rho_i|a_i)$  stands for the probability that  $i$  deviates to  $(b_i, \rho_i)$  when recommended to play  $a_i$ . Given  $\sigma \in \Delta(A)$ , let  $\Pr(\sigma) \in \mathbb{R}^S$  be the vector defined by  $\Pr(\sigma)(s) = \sum_a \sigma(a) \Pr(s|a)$ . Intuitively,  $\Pr(\sigma)$  is the vector of prior report probabilities if everyone is honest and obediently playing according to  $\sigma$ . Let  $\Pr(\sigma, \alpha_i) \in \mathbb{R}^S$ , defined pointwise by

$$\Pr(\sigma, \alpha_i)(s) = \sum_{a \in A} \sigma(a) \sum_{(b_i, \rho_i)} \Pr(s|a_{-i}, b_i, \rho_i) \alpha_i(b_i, \rho_i|a_i),$$

be the vector of prior probabilities if agent  $i$  deviates from  $\sigma$  according to  $\alpha_i$ .

A deviation plan  $\alpha_i$  is *disobedient* if  $\alpha_i(b_i, \rho_i|a_i) > 0$  for some  $a_i \neq b_i$ , i.e., it disobeys some recommendation  $a_i$  with positive probability. A disobedient deviation plan may be “honest,” i.e.,  $\rho_i$  may equal  $\tau_i$  with probability one after every recommendation. A profile  $\alpha = (\alpha_1, \dots, \alpha_n)$  of deviation plans is called *disobedient* if  $\alpha_i$  is disobedient for some agent  $i$ . Although dishonesty is arguably a form of disobedience, it will be useful in the sequel to distinguish between them.

**Definition 2 (Detection).** A deviation plan  $\alpha_i$  for agent  $i$  is called *undetectable* if

$$\forall \sigma \in \Delta(A), \quad \Pr(\sigma) = \Pr(\sigma, \alpha_i).$$

Call  $\alpha_i$  *detectable* if it is not undetectable, i.e.,  $\Pr(\sigma) \neq \Pr(\sigma, \alpha_i)$  for some  $\sigma \in \Delta(A)$ .

Intuitively, a deviation plan  $\alpha_i$  is undetectable if the probability of reported signals induced by  $\alpha_i$ ,  $\Pr(\sigma, \alpha_i)$ , coincides with that arising from honesty and obedience,  $\Pr(\sigma)$ , *regardless of the team’s correlated strategy*,  $\sigma$ , assuming that others are honest and obedient. Undetectability is arguably a strong restriction on a deviation plan, making detectability a weak requirement.<sup>10</sup> We now give our first main definition.

**Definition 3 (DUD).** A monitoring technology  $\Pr$  *detects unilateral disobedience* (DUD) if every disobedient deviation plan is detectable.

DUD is intuitively defined.<sup>11</sup> Formally, note that different correlated strategies may be used to decide whether or not different disobedient deviation plans are detectable.

<sup>10</sup>Undetectability may be defined equivalently by  $\Pr(a) = \Pr(a, \alpha_i)$  for all  $a \in A$  by linearity.

<sup>11</sup>For a slightly stronger but also mathematically more tractable version of DUD (without using reporting strategies), see Lemma B.1.

This is one important aspect that renders DUD substantially weaker than other conditions in the literature, as will soon be seen. To illustrate, consider an example.

**Example 1.** There are two publicly verifiable signals,  $S = S_0 = \{x, y\}$ , and two agents,  $I = \{1, 2\}$ . Agent 1 has two actions,  $A_1 = \{U, D\}$ , and agent 2 has three actions,  $A_2 = \{L, M, R\}$ . The conditional probability system  $\Pr$  is given below.

	$L$	$M$	$R$
$U$	1, 0	0, 1	1/2, 1/2
$D$	1, 0	0, 1	1/3, 2/3

If agent 1 plays  $U$  then there is a mixed deviation by agent 2 (namely  $\frac{1}{2}[L] + \frac{1}{2}[M]$ , where  $[\cdot]$  stands for Dirac measure) such that the conditional probability over signals equals what it would be if he played  $R$ . A similar phenomenon takes place when agent 1 plays  $D$  (this time with the deviation  $\frac{2}{3}[L] + \frac{1}{3}[M]$ ) or indeed regardless of agent 1's mixed strategy. It is therefore impossible to even approximately enforce  $R$  with transfers contingent only on signals if agent 2 strictly prefers playing  $L$  and  $M$ , since there always exists a profitable deviation without any contractual losses.

However,  $\Pr$  detects unilateral disobedience because for any deviation plan by agent 2 there is a mixed strategy by agent 1 that detects it. By correlating agent 2's payment with agent 1's recommendation, secret contracts can keep agent 2 from knowing the proportion with which he ought to mix between  $L$  and  $M$  for his contractual payment to equal what he would obtain by playing  $R$ . It will be seen that this renders  $R$  enforceable. This suggests how secret contracts can extract more information from a monitoring technology to provide incentives, even with publicly verifiable signals.

Next, we will show that DUD characterizes approximate enforcement.

**Definition 4 (PSI).** A monitoring technology  $\Pr$  *provides strict incentives* (PSI) if there exists a probability weighted incentive scheme  $\xi : I \times A \times S \rightarrow \mathbb{R}$  such that

$$\forall(i, a_i, b_i, \rho_i), \quad 0 \leq \sum_{(a_{-i}, s)} \xi_i(a, s)(\Pr(s|a_{-i}, b_i, \rho_i) - \Pr(s|a)),$$

with a *strict inequality* whenever  $a_i \neq b_i$ .<sup>12</sup>

---

<sup>12</sup>Although no budget constraints are imposed, we could have added *expected* budget balance,

$$\sum_{(i, a, s)} \xi_i(a, s) = 0,$$

but this constraint would not bind, since adding a constant to any  $\xi$  preserves its incentive properties.

By scaling  $\xi$  as necessary, PSI implies that for every utility profile there is an incentive scheme so that any deviator's contractual loss outweighs his deviation gain. PSI may appear to be a rather strong condition, in contrast with the argued weakness of DUD (Example 1). As it turns out, PSI and DUD are equivalent, in fact mutually dual.

**Lemma 1.** *A monitoring technology detects unilateral disobedience if and only if it provides strict incentives.*

*Proof.* By the Alternative Theorem (Rockafellar, 1970), PSI fails if and only if there is a vector  $\lambda \geq 0$  such that  $\lambda_i(a_i, b_i, \rho_i) > 0$  for some  $(i, a_i, b_i, \rho_i)$  with  $a_i \neq b_i$  and

$$\forall(a, s), \quad \sum_{(b_i, \rho_i)} \lambda_i(a_i, b_i, \rho_i) (\Pr(s|a_{-i}, b_i, \rho_i) - \Pr(s|a)) = 0.$$

Such a vector  $\lambda$  exists if and only if the deviation plan  $\alpha_i$ , defined pointwise by

$$\alpha_i(b_i, \rho_i|a_i) := \begin{cases} \lambda_i(a_i, b_i, \rho_i) / \sum_{(b'_i, \rho'_i)} \lambda_i(a_i, b'_i, \rho'_i) & \text{if } \sum_{(b'_i, \rho'_i)} \lambda_i(a_i, b'_i, \rho'_i) > 0, \text{ and} \\ [(a_i, \tau_i)](b_i, \rho_i) & \text{otherwise (where } [\cdot] \text{ denotes Dirac measure),} \end{cases}$$

is disobedient and undetectable: DUD fails.  $\square$

The simple proof of Lemma 1 above describes a duality between identifiability and enforceability via secret contracts. A natural corollary follows that motivates DUD from a “backward-engineering” exercise: what minimal requirement on a monitoring technology suffices to contractually overcome incentive constraints? Given  $\xi$  and any completely mixed correlated strategy  $\sigma \in \Delta^0(A) := \{\sigma \in \Delta(A) : \forall a \in A, \sigma(a) > 0\}$ , there exists  $\zeta$  with  $\xi_i(a, s) = \sigma(a)\zeta_i(a, s)$  for all  $(i, a, s)$ . Hence, PSI is equivalent to every  $\sigma \in \Delta^0(A)$  being (exactly) enforceable, which proves the next result.

**Theorem 1.** *A monitoring technology detects unilateral disobedience if and only if any team with any profile of utility functions can approximately enforce any correlated strategy with secret contracts.*

As Example 1 shows, DUD is not enough to provide incentives with just signal-contingent contracts, but the following strengthening is. Given a subset  $B \subset A$  and an agent  $i$ , let  $B_i := \{b_i \in A_i : \exists b_{-i} \in A_{-i} \text{ s.t. } b \in B\}$  be the projection of  $B$  on  $A_i$ . Call a deviation plan  $\alpha_i$  *B-disobedient* if it is disobedient at some  $a_i \in B_i$ . Given  $\sigma \in \Delta(A)$ , say  $\Pr$  *detects unilateral disobedience* at  $\sigma$  (DUD- $\sigma$ ) if  $\Pr(\sigma) \neq \Pr(\sigma, \alpha_i)$  for every agent  $i$  and  $\text{supp } \sigma$ -disobedient<sup>13</sup> deviation plan  $\alpha_i$ . Intuitively, the same  $\sigma$  detects every  $\alpha_i$ . The proof of Theorem 1 also implies the following corollary.

<sup>13</sup>By definition,  $\text{supp } \sigma = \{a \in A : \sigma(a) > 0\}$  is the support of  $\sigma$ .

**Corollary 1.** *Fix a correlated strategy  $\sigma$ . A monitoring technology detects unilateral disobedience at  $\sigma$  if and only if any team with any profile of utility functions can enforce  $\sigma$  with just “standard” signal-contingent contracts.*

Corollary 1 captures the value-added of secret contracts. By the paragraph preceding Theorem 1, DUD suffices to enforce any completely mixed  $\sigma$  with secret transfers by effectively allowing the use of different  $\sigma$ ’s to detect different  $\alpha$ ’s, unlike standard contracts, for which the same  $\sigma$  must detect every  $\alpha$ .<sup>14</sup> On the other hand, to enforce a pure-strategy profile  $a$ , DUD is generally not enough. Since agents receive only one recommendation under  $[a]$ , there is no use for secret contracts, so by Corollary 1 DUD- $[a]$  characterizes enforcement with secret as well as standard contracts.

The intermediate case where  $\sigma$  has arbitrary support is discussed in Section 4.1. There, necessary and sufficient conditions are derived for exact as well as approximate enforcement. Section 4.2 extends the results further by fixing utility functions.

Genericity of DUD is established next from the number of agents’ action-signal pairs. Intuitively, incentives may be provided to a given agent in three ways: (a) using only others’ signals to detect his deviations (e.g., Friday), (b) using only his own reports and others’ recommendations (e.g., Robinson), and (c) using both his reports and others’ signals in conjunction. Theorem 2 below identifies conditions such that for every agent, at least one such way of detecting deviations is generic.

**Theorem 2.** *DUD is generic if for every agent  $i$ ,*

- (a)  $|A_i| - 1 \leq |A_{-i}| (|S_{-i}| - 1)$  when  $|S_i| = 1$ ,
- (b)  $|A_i| (|S_i| - 1) \leq |A_{-i}| - 1$  when  $|S_{-i}| = 1$ , and
- (c)  $|A_i| |S_i| \leq |A_{-i}| |S_{-i}|$  when both  $|S_i| > 1$  and  $|S_{-i}| > 1$ .

If  $|S| = 1$  then DUD is generic only if  $|A| = 1$ . More interestingly, DUD is generic even if  $|S| = 2$ , as long as agents have enough actions. Hence, a team may overcome incentive constraints (i.e., DUD, therefore Theorem 1, holds) generically even if only one individual can make substantive observations and these observations are just a binary bit of information. If others’ action spaces are large enough and their actions have generic effect on the bit’s probability, this uniquely informed individual may still be controlled by testing him with unpredictable combinations of others’ actions.<sup>15</sup>

---

<sup>14</sup>Even for approximate enforcement with standard contracts the same  $\sigma$  must detect all  $\alpha$ ’s. E.g., in Example 1 there is no sequence  $\{\sigma^m\}$  with  $\sigma^m \rightarrow [(U, R)]$  and  $\Pr$  satisfying DUD- $\sigma^m$  for all  $m$ .

<sup>15</sup>We thank an anonymous referee for urging us to emphasize this point.

We conclude this subsection by relating DUD to the literature. Broadly, DUD is an improvement in that different  $\sigma$  can be used to detect different  $\alpha_i$ .

In a restricted setting, Legros and Matsushima (1991) and Legros and Matthews (1993) find conditions equivalent to DUD-[ $a$ ] (but differently interpreted) to enforce a profile  $a$  with signal-contingent contracts. In repeated games, Fudenberg et al. (1994) introduced *individual full rank* (IFR). Formally, IFR (at some  $\sigma$ ) means that for every  $i$ ,  $\Pr(\sigma) \notin \text{span}\{\Pr(\sigma, b_i, \rho_i) : (b_i, \rho_i) \neq (a_i, \tau_i)\}$ , where “span” stands for linear span. Arguably, the spirit of IFR is to detect deviations away from some prescribed  $\sigma$ , i.e., DUD- $\sigma$ .<sup>16</sup> IFR at  $\sigma$  implies DUD- $\sigma$  but not conversely.<sup>17</sup> DUD is also weaker than *local* IFR (LIFR) of d’Aspremont and Gérard-Varet (1998), requiring IFR at possibly different  $\sigma$  for different  $i$ .<sup>18</sup> Indeed, clearly LIFR implies DUD, and LIFR fails but DUD holds in Example 1. “Local” DUD- $\sigma$  fails there, too.

DUD is also weaker than the generalization of IFR by Kandori (2003), where agents play mixed strategies and report on the realization of such mixtures. He considers contracts contingent on those reports and signal realizations. The next example shows that secret contracts can perform strictly better even with public monitoring.

**Example 2.** One agent, three actions ( $L$ ,  $M$  and  $R$ ), and two publicly verifiable signals ( $g$  and  $b$ ), with the following utility function and monitoring technology.

$L$	$M$	$R$	$L$	$M$	$R$
0	2	0	1, 0	1/2, 1/2	0, 1

Utility Payoffs                      Signal Probabilities

The mixed strategy  $\sigma = \frac{1}{2}[L] + \frac{1}{2}[R]$  is enforceable with secret contracts but not with Kandori’s contracts. Indeed, offering \$1 for  $g$  if asking to play  $L$  and \$1 for  $b$  if asking to play  $R$  makes  $\sigma$  enforceable. With Kandori’s contracts, the agent is asked to play  $\sigma$  and then asked what he actually played before receiving any monetary rewards. The agent gains two ‘utils’ by playing  $M$  instead and announcing that he played  $L$  ( $R$ ) if the realized signal is  $g$  ( $b$ ), with the same expected monetary payoff.<sup>19</sup>

<sup>16</sup>For instance, see Compte (1998) or Kandori and Matsushima (1998).

<sup>17</sup>If  $|S_{-i}| < |A_i|$  for some  $i$  then this holds trivially, since IFR is impossible yet DUD- $\sigma$ , which requires only convex (rather than linear) independence, is possible (e.g., all the points on a circle are convexly independent). This holds even with at least as many signals as actions (e.g., consider the vectors  $(\frac{1}{3}, \frac{1}{3}, 0, \frac{1}{3})$ ,  $(0, \frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ ,  $(\frac{1}{6}, 0, \frac{1}{3}, \frac{1}{2})$  and  $(\frac{1}{3}, 0, \frac{1}{6}, \frac{1}{2})$ ).

<sup>18</sup>For all  $i$ , LIFR uses the same correlated strategy  $\sigma_i$  to detect each deviation plan  $\alpha_i$  of agent  $i$ .

<sup>19</sup>Kandori’s are like secret contracts if actions are *secretly* announced *before* signals are observed.

Tomala (2005) independently derives a condition comparable to DUD to prove a folk theorem. He defines detection with respect to a fixed correlated strategy using unconditional probabilities over actions and signals. He focuses on exact implementation, so for  $\sigma \in \Delta^0(A)$ , his version of DUD agrees broadly with ours (he proves a version of Corollary 3 below). However, he does not study approximate enforcement in general (i.e., for  $\sigma \notin \Delta^0(A)$ ), and does not use different  $\sigma$  to detect different  $\alpha_i$ . This issue is developed further in Theorems 3 and 5 (Section 4) below.

## 4 Discussion

This section makes two comments. Firstly, the previous section's results are extended to correlated strategies with restricted support. Finally, we discuss collusion and characterize contracts that dissuade multilateral deviations.

### 4.1 Exact versus Approximate Enforcement

Next, we characterize exact enforcement of a fixed correlated strategy for any utility profile. Fix two sets of action profiles  $B, C \subset A$ . A deviation plan  $\alpha_i$  is *C-detectable* if  $\Pr(\sigma) \neq \Pr(\sigma, \alpha_i)$  for some  $\sigma \in \Delta(A)$  with  $\text{supp } \sigma \subset C$ . Say  $\Pr$  *C-detects unilateral B-disobedience* ( $\text{DUD}_B^C$ ) if every *B-disobedient* deviation plan is *C-detectable*. (We will call *A-detection* simply *detection*, and write  $\text{DUD}_B^A$  as  $\text{DUD}_B$ .) For instance,  $\text{DUD}_A^A$  is just DUD, and  $\text{DUD}_{\{a\}}^{\{a\}}$  equals  $\text{DUD}-[a]$ . Consider another example.

**Example 3.** There are two agents and two publicly verifiable signals, with the monitoring technology below. (It is Example 1 with an added row.)

	$L$	$M'$	$R$
$U$	1, 0	0, 1	2/3, 1/3
$M$	1, 0	0, 1	1/2, 1/2
$D$	1, 0	0, 1	1/3, 2/3

Let  $A = \{U, M, D\} \times \{L, M', R\}$ ,  $B = \{U, M\} \times \{L, M', R\}$ ,  $C = \{U, D\} \times \{L, M', R\}$ . Clearly,  $\text{DUD}_A^A$  fails here, since  $\frac{1}{2}[U] + \frac{1}{2}[D]$  is statistically indistinguishable from  $M$ . Also,  $\text{DUD}_B^A$  fails because a plan to play  $\frac{1}{2}[U] + \frac{1}{2}[D]$  when asked to play  $M$  is *A-undetectable* and *B-disobedient*. However,  $\text{DUD}_C^C$  does hold, since there is no undetectable deviation from  $D$  or  $U$  by the row player. ( $\text{DUD}-\sigma$  fails for every  $\sigma$ .)

**Corollary 2.** *Fix any subset  $B \subset A$ . A monitoring technology  $B$ -detects unilateral  $B$ -disobedience if and only if any team with any profile of utility functions can exactly enforce every (if and only if some fixed) correlated strategy with support equal to  $B$ .*

Therefore, every correlated strategy with support equal to a subset of  $B$  is automatically approximately enforceable, just as with Theorem 1. By Corollary 2, *only the support of a correlated strategy matters* for its enforcement regardless of preferences.

Having characterized exact enforcement, we proceed with approximate enforcement. By Corollary 2, existence of some  $C \supset B$  such that  $\text{DUD}_C^C$  holds clearly yields a sufficient condition. However, this is unnecessary. To motivate, consider an example.

**Example 4.** Two agents, two public signals, the following monitoring technology:

	$L$	$M$	$R$
$U$	1, 0	1, 0	1, 0
$D$	1, 0	0, 1	0, 1

The action profile  $(U, L)$  is not enforceable for every utility profile, since  $\text{DUD}_{\{(U,L)\}}^{\{(U,L)\}}$  clearly fails. Indeed, playing  $D$  when asked to play  $U$  is  $\{(U, L)\}$ -undetectable. It is also easy to see that  $\text{DUD}_C^C$  fails, too, for every  $C \supset \{(U, L)\}$ . However,  $(U, L)$  is approximately enforceable for every utility profile, since either  $[(D, M)]$  or  $[(D, R)]$  can be used to detect  $\{(U, L)\}$ -disobedient deviations. No contract can induce the row player to choose  $M$  if  $R$  is strictly dominant, say, but this is unimportant as long as the row player chooses either  $M$  or  $R$  when asked to do so. The key condition satisfied here is that every  $\{(U, L)\}$ -disobedient deviation plan is detectable ( $\text{DUD}_{\{(U,L)\}}$ ).

In general, it should be clear that  $\text{DUD}_B$  is necessary for approximate enforcement, but perhaps it is not so clear that it is also sufficient, as the next result shows.

**Theorem 3.** *Fix any subset  $B \subset A$ . A monitoring technology detects unilateral  $B$ -disobedience if and only if any team with any profile of utility functions can approximately enforce every (if and only if some fixed) correlated strategy with support in (equal to)  $B$  with secret contracts.*

Hence, still only the support of a correlated strategy matters for approximate enforcement regardless of preferences. Clearly, Theorem 1 is a special case of Theorem 3 (as well as Corollary 2) when  $B = A$ . Example 4 illustrates the insight behind Theorem 3 and gives intuition for its proof. Suppose that, to detect deviations from  $a_i \in B_i$ ,



some  $a_j \notin B_j$  is played infrequently by  $j \neq i$ . What if  $a_j$  itself has a profitable, undetectable deviation  $\alpha_j(a_j) \in \Delta(A_j)$ ? After all,  $\text{DUD}_B$  says nothing about detection outside  $B$ . If such  $\alpha_j(a_j)$  exists, playing it instead of  $a_j$  *still detects deviations* from  $a_i$  by virtue of being undetectable. Similarly, undetectable deviations from  $\alpha_j(a_j)$  detect deviations from  $a_i$ , and so on. Proceeding iteratively, since the game is finite there must be detecting behavior without a profitable, undetectable deviation.

## 4.2 Fixed Utility Functions

Throughout this section, let us fix a profile  $v : I \times A \rightarrow \mathbb{R}$  of utility functions. A natural weakening of the previous results might be to allow for undetectable deviation plans as long as they are *unprofitable*. Exact enforcement amounts to the following.

**Corollary 3.** *A correlated strategy  $\sigma$  is enforceable with secret contracts if and only if every supp  $\sigma$ -undetectable deviation plan  $\alpha_i$  of any agent  $i$  is  $\sigma$ -unprofitable, i.e.,*

$$\Delta v_i(\sigma, \alpha_i) := \sum_{(a, b_i, \rho_i)} \sigma(a) \alpha_i(b_i, \rho_i | a_i) (v_i(a_{-i}, b_i) - v_i(a)) \leq 0.$$

The proof of this claim is comparable to previous ones, therefore omitted. Given an enforceable correlated strategy  $\sigma$ , we now ask how large transfers must be to enforce it. To this end, let us introduce some notation. Let  $\mathcal{D}_i = \Delta(A_i \times R_i)^{A_i}$  be the space of deviation plans  $\alpha_i$  for a agent  $i$  and  $\mathcal{D} = \prod_i \mathcal{D}_i$  be the set of profiles of deviation plans  $\alpha = (\alpha_1, \dots, \alpha_n)$ . For any deviation plan  $\alpha_i \in \mathcal{D}_i$  and any  $\sigma \in \Delta(A)$ , write

$$\|\Delta \text{Pr}(\sigma, \alpha_i)\| := \sum_{s \in S} \left| \sum_{(a, b_i, \rho_i)} \sigma(a) (\alpha_i(b_i, \rho_i | a_i) \text{Pr}(s | a_{-i}, b_i, \rho_i) - \text{Pr}(s | a)) \right|.$$

This norm summarizes the difference in signal probabilities between abiding by  $\sigma$  and deviating to  $\alpha_i$ . A correlated strategy  $\sigma$  is called *enforceable within* some vector  $z \in \mathbb{R}_+^I$  if there exists a scheme  $\xi : I \times A \times S \rightarrow \mathbb{R}$  to satisfy  $(*)$  and

$$\forall(i, a, s), \quad -\sigma(a)z_i \leq \xi_i(a, s) \leq \sigma(a)z_i.$$

Next, we provide a lower bound on  $z$  so that  $\sigma$  is enforceable within  $z$ .

**Theorem 4.** *(i) A correlated strategy  $\sigma$  is enforceable within  $z \in \mathbb{R}_+^I$  if and only if*

$$V_\sigma(z) := \max_{\alpha \in \mathcal{D}} \sum_{i \in I} \Delta v_i(\sigma, \alpha_i) - \sum_{(i, a)} z_i \sigma(a) \|\Delta \text{Pr}(a, \alpha_i)\| = 0.$$

- (ii) If  $\sigma$  is enforceable then  $V_\sigma(z) = 0$  for some  $z \in \mathbb{R}_+^I$ . If not then  $\sup_z V_\sigma(z) > 0$ .  
(iii) A correlated strategy  $\sigma$  is enforceable if and only if  $\bar{z}_i < +\infty$  for every  $i$ , where

$$\bar{z}_i := \sup_{\alpha_i \in \mathcal{F}_i} \frac{\max\{\Delta v_i(\sigma, \alpha_i), 0\}}{\sum_a \sigma(a) \|\Delta \Pr(a, \alpha_i)\|} \quad \text{if } \mathcal{F}_i := \{\alpha_i : \sum_a \sigma(a) \|\Delta \Pr(a, \alpha_i)\| > 0\} \neq \emptyset$$

and, whenever  $\mathcal{F}_i = \emptyset$ ,  $\bar{z}_i := +\infty$  exactly when  $\max_{\alpha_i} \Delta v_i(\sigma, \alpha_i) > 0$ .<sup>20</sup>

- (iv) If  $\bar{z}_i < +\infty$  for every  $i$  then  $V_\sigma(z) = 0$  if and only if  $z_i \geq \bar{z}_i$  for all  $i$ .

Theorem 4 quantifies the wedge that transfers require to enforce a given correlated strategy with punishments and rewards. It implies that  $\text{supp } \sigma$ -detectability, hence also enforceability by Corollary 3, is captured by  $\sum_a \sigma(a) \|\Delta \Pr(a, \alpha_i)\| > 0$ . By Corollary 1, enforcement with signal-contingent transfers is captured by the stronger detectability condition that  $0 < \|\Delta \Pr(\sigma, \alpha_i)\| \leq \sum_a \sigma(a) \|\Delta \Pr(a, \alpha_i)\|$ . Hence, a version of Theorem 4 holds with signal-contingent transfers and  $\|\Delta \Pr(\sigma, \alpha_i)\|$  instead of  $\sum_a \sigma(a) \|\Delta \Pr(a, \alpha_i)\|$ . Finally, Theorem 4 (iii) clearly implies the following.

**Corollary 4.** *Each  $\text{supp } \sigma$ -undetectable deviation plan is  $\sigma$ -unprofitable if and only if there exists  $z \geq 0$  such that  $\Delta v_i(\sigma, \alpha_i) \leq z \sum_a \sigma(a) \|\Delta \Pr(a, \alpha_i)\|$  for all  $i$  and  $\alpha_i$ , that is, utility gains from any deviation are uniformly outweighed by monetary losses.*

Next, we characterize approximate enforcement of a correlated strategy for a fixed profile of utility functions. Corollary 3 might suggest that a correlated strategy  $\sigma$  is enforceable if every  $\sigma$ -profitable deviation plan is detectable. However, the next example shows that approximate enforcement imposes a stronger requirement.

**Example 5.** Consider a variation on Robinson and Friday (Section 2.1):

	work	shirk	solitaire		work	shirk	solitaire
monitor	0, -1	0, 0	0, 0	monitor	1, 0	0, 1	1, 0
rest	0, -1	0, 0	0, -1	rest	1/2, 1/2	1/2, 1/2	1/2, 1/2
Utility Payoffs				Signal Probabilities			

Assume for simplicity that the signal is *publicly verifiable* and Robinson's utility is constant. Clearly, the profile (rest, work) is not enforceable because a deviation by Friday to shirk is [(rest, work)]-profitable and {(rest, work)}-undetectable. Moreover, (rest, work) is *not approximately enforceable* either. Indeed, for Friday to ever

<sup>20</sup>Intuitively,  $\mathcal{F}_i$  is the set of all  $\text{supp } \sigma$ -detectable deviation plans available to agent  $i$ .

work Robinson must monitor with positive probability. But then no contract can discourage Friday from playing solitaire instead of working, since playing solitaire when asked to work is undetectable and weakly dominant. On the other hand, every  $[(\text{rest}, \text{work})]$ -profitable deviation plan is detectable.

Removing solitaire from [Example 5](#) restores approximate enforcement of  $(\text{rest}, \text{work})$ . This occurs not because every  $(\text{rest}, \text{work})$ -profitable deviation is detectable (it is true with or without solitaire), but because it is *uniformly detectable*, i.e., the utility gains from every  $(\text{rest}, \text{work})$ -profitable deviation by Friday are uniformly outweighed by monetary losses when Robinson monitors, in line with [Corollary 4](#). The next result characterizes approximate enforcement with “uniform, credible” detection.

**Theorem 5.** *A correlated strategy  $\sigma$  is approximately enforceable if and only if there exists  $z \geq 0$  such that every  $\sigma$ -profitable deviation plan  $\alpha_i$  is detectable by some correlated strategy  $\mu$  for which both*

- (i)  $\Delta v_i(\mu, \alpha_i) < z \sum_a \mu(a) \|\Delta \Pr(a, \alpha_i)\|$  and
- (ii)  $\Delta v_j(\mu, \alpha_j) \leq z \sum_a \mu(a) \|\Delta \Pr(a, \alpha_j)\|$  for every other agent  $j$  and deviation  $\alpha_j$ .

Intuitively, [Theorem 5](#) says that to approximately enforce a correlated strategy, it is necessary and sufficient that all its profitable deviations be discouraged both (i) uniformly and (ii) credibly. As before, different behavior may be used to detect different deviations by an agent.<sup>21</sup> Formally, uniform detection means that for the same fixed  $z$ , every deviation plan  $\alpha_i$  must impact the magnitude of  $z$ -weighted probabilistic changes enough to outweigh its deviation gains. Therefore, transfers bounded within  $z$  can provide incentives against all  $\sigma$ -profitable deviations, perhaps with different  $\mu$  for different  $\alpha_i$ .

To explain the need for credibility, compare this result with [Theorem 3](#), where “credible monitoring” is unnecessary. There, *every* disobedient deviation is potentially profitable, so ought to be detectable. Here, with fixed utility functions, even if some disobedient deviation plan  $\alpha_i$  is undetectable, it may nonetheless be discouraged with behavior  $\mu$  by others that makes the deviation unprofitable (as in a correlated equilibrium without transfers). However, if this specific behavior is not credible then there may exist a  $\mu$ -profitable deviation plan  $\alpha_j$  by some other agent such that  $\alpha_i$  becomes

---

<sup>21</sup>To see that credibility matters, simply add a row to the table in [Example 5](#) above with utility payoffs 

$-1, -1$	$-1, 0$	$-1, -1$
----------	---------	----------

 and signal probabilities 

$1, 0$	$0, 1$	$1, 0$
--------	--------	--------

. Now there is an action for Robinson that is strictly dominated and indistinguishable from monitoring, yet uniformly detects all of Friday’s  $(\text{rest}, \text{work})$ -profitable deviations.

profitable once again given  $\mu$  and  $\alpha_j$ .

We end this subsection by noting without proof that all previous results hold also with budget balance using the same arguments, replacing detection with attribution and  $\sum_a \sigma(a) \|\Delta \Pr(a, \alpha_i)\|$  with  $\min_{\eta} \sum_{(i,a)} \sigma(a) \|\Delta \Pr(a, \alpha_i) - \eta(a)\|$ , where  $\eta \in \mathbb{R}^{A \times S}$ . Now, instead of the change in transfers to an agent after a deviation being bounded by the magnitude of the change in the probability over signals, it is bounded by the residuals of a least-absolute-deviations regression of the probability changes on  $A$ . This amount is clearly smaller than the magnitude of the dependent variable, i.e., the probability change. Intuitively, budget balance implies that only the variation across agents of the deviations' effect on signals can be used to provide incentives to discourage them, rather than the deviations' effects themselves.

### 4.3 Coalitional Deviations

A notable weakness of secret contracts is not being collusion-proof. To illustrate, in our leading example (Section 2.1) Robinson and Friday could communicate “extra-contractually” to break down the incentives that secrets tried to provide.<sup>22</sup> On the other hand, collusion is a problem for contracts in general. For instance, the scheme proposed by Cremer and McLean (1988) is not collusion-proof for similar reasons.

To study collusion-proof contracts, assumptions must be made regarding coalitions' contractual ability. We will assume that every coalition  $t$  maximizes some given *coalitional utility function*  $v_t : A \rightarrow \mathbb{R}$ , quasilinear in monetary transfers.<sup>23</sup>

**Definition 5.** A correlated strategy  $\sigma$  is *strongly enforceable* if there is an incentive scheme  $\zeta : I \times A \times S \rightarrow \mathbb{R}$  such that

$$\forall t \subset I, a_t \in A_t, (b_t, \rho_t) \in A_t \times R_t, \\ \sum_{a-t} \sigma(a) (v_t(a_{-t}, b_t) - v_t(a)) \leq \sum_{(a-t, s)} \sigma(a) \sum_{i \in t} \zeta_i(a, s) (\Pr(s|a_{-t}, b_t, \rho_t) - \Pr(s|a)).$$

<sup>22</sup>The following incentive scheme deters such communication between Robinson and Friday (Friday prefers misreporting his signal to Robinson) while approximately enforcing (rest, work).

	(monitor, work)	(monitor, shirk)	(rest, work)	(rest, shirk)
$g$	$1/\mu, 1/\sigma$	$0, 1/\sigma$	$1/2\mu, 0$	$0, 1/2(1 - \sigma)$
$b$	$0, 0$	$1/(1 - \mu), 0$	$0, 1/(1 - \sigma)$	$1/2(1 - \mu), 1/2(1 - \sigma)$

<sup>23</sup>This assumption is standard. See for instance, Che and Kim (2006) and references therein. The purpose of this section is not to derive a meaningful utility for coalitions, but to use one.

Strong enforcement requires that no subset of agents can profitably deviate after coordinating their information even if they can commit to sharing their information non-strategically. This makes strong enforceability especially “strong.”

We now derive the detection requirement implied by strong enforceability. Given a nonempty subset of agents  $t \subset I$ , a *multilateral deviation plan* for  $t$  is any measure-valued map  $\alpha_t : A_t \rightarrow \Delta(A_t \times R_t)$ , where  $A_t \times R_t = \prod_{i \in t} A_i \times R_i$ . Intuitively, a multilateral deviation plan  $\alpha_t$  has the agents in  $t$  coordinate their deviations contingent on all recommendations to members of  $t$ . A multilateral deviation plan  $\alpha_t$  is called *disobedient* if  $\alpha_t(b_t, \rho_t | a_t) > 0$  for some  $(a_t, b_t, \rho_t)$  such that  $a_t \neq b_t$ . It is called *detectable* if  $\Pr(\sigma) \neq \Pr(\sigma, \alpha_t)$  for some  $\sigma \in \Delta(A)$ .

A *coalitional deviation plan* by agent  $i$  is a profile of multilateral deviation plans  $\alpha^i = \{\alpha_t : t \ni i\}$ , one for each coalition to which  $i$  may belong. It is called *disobedient* if  $\alpha_t$  is disobedient for some coalition  $t \ni i$ . It is called *detectable* if  $\Pr(\sigma) \neq \Pr(\sigma, \alpha^i)$  for some  $\sigma \in \Delta(A)$ , where  $\Pr(\sigma, \alpha^i) := \sum_{t \ni i} \sum_{(a, b_t, \rho_t)} \sigma(a) \Pr(a_{-t}, b_t, \rho_t) \alpha_t(b_t, \rho_t | a_t)$ . Intuitively, a coalitional deviation plan for an agent  $i$  is a profile of multilateral deviation plans involving  $i$ . It is undetectable if regardless of the correlated strategy  $\sigma$ , even if some multilateral deviation plan  $\alpha_t$  is detectable, there is another multilateral deviation plan  $\alpha_{t'}$  with  $i \in t \cap t'$  that “undoes” the change in probability from  $\alpha_t$ . Therefore, even if every disobedient multilateral deviation plan is detectable, it is possible that some disobedient coalitional deviation plan remains undetectable.

**Definition 6 (DCD).** A monitoring technology  $\Pr$  *detects coalitional disobedience* (DCD) if every disobedient coalitional deviation plan is detectable.

The next result characterizes strong enforcement as detection of coalitional deviations. It is argued similarly to previous ones, so its proof is omitted.

**Theorem 6.** *A monitoring technology detects coalitional disobedience if and only if any team with any profile of coalitional utility functions can approximately strongly enforce every correlated strategy with secret contracts.*

## 5 Conclusion

Secret contracts emphasize that—as part of a team’s economic organization—it may be beneficial for private information to be allocated differently across individuals in

order for the right incentives to be provided. This remains true even if the team starts without informational asymmetry. Secret contracts effectively subject contractual deviations to “tailored monitoring,” making monitors de facto auditors. Heuristically, secret contracts allow for monitoring to *follow* deviations.

## A Proofs

*Corollary 1.* Fix any  $\sigma \in \Delta(A)$ . By the Alternative Theorem, a monitoring technology  $\Pr$  satisfies DUD- $\sigma$  if and only if there is a signal-contingent scheme  $\zeta : I \times S \rightarrow \mathbb{R}$  such that

$$\forall i \in I, a_i \in B_i, (b_i, \rho_i) \in A_i \times R_i, \quad 0 \leq \sum_{(a_{-i}, s)} \sigma(a) \zeta_i(s) (\Pr(s|a_{-i}, b_i, \rho_i) - \Pr(s|a)),$$

with a strict inequality if  $a_i \neq b_i$ , where  $B_i = \{a_i \in A_i : \exists a_{-i} \text{ s.t. } \sigma(a) > 0\}$ . Call this dual condition PSI- $\sigma$ . By scaling  $\zeta$  appropriately, PSI- $\sigma$  clearly implies that any deviation gains can be outweighed by contractual losses. Conversely, if DUD- $\sigma$  fails then  $\Pr(\sigma) = \Pr(\sigma, \alpha_i)$  for some deviation plan  $\alpha_i$  with  $\alpha_i(b_i, \rho_i|a_i) > 0$  for some  $a_i \in B_i$ , and  $b_i \neq a_i$ . For all  $a_{-i}$ , let  $0 = v_i(a) < v_i(a_{-i}, b_i) = 1$ . Now  $\sigma$  cannot be enforced by any  $\zeta : I \times S \rightarrow \mathbb{R}$ , since  $\sum_{(b_i, \rho_i)} \alpha_i(b_i, \rho_i|a_i) \sum_{a_{-i}} \sigma(a) (v_i(a_{-i}, b_i) - v_i(a)) > \sum_s \zeta_i(s) (\Pr(s|\sigma, \alpha_i) - \Pr(s|\sigma)) = 0$ , being a convex combination of incentive constraints, must violate at least one.  $\square$

*Theorem 2.* By Lemma B.1, DUD is implied by *conic independence*

$$\forall (i, a_i, s_i), \quad \Pr(a_i, s_i) \notin \text{cone}\{\Pr(b_i, t_i) : (b_i, t_i) \neq (a_i, s_i)\}.$$

This is in turn implied by *linear independence*, or full row rank, for all  $i$ , of the matrix with  $|A_i| |S_i|$  rows,  $|A_{-i}| |S_{-i}|$  columns and entries  $\Pr(a_i, s_i)(a_{-i}, s_{-i}) = \Pr(s|a)$ . Since the set of full rank matrices is generic, this full row rank is generic if  $|A_i| |S_i| \leq |A_{-i}| |S_{-i}|$  if  $|S_i| > 1$  and  $|S_{-i}| > 1$ . If  $|S_i| = 1$ , adding with respect to  $s_{-i}$  for each  $a_{-i}$  yields column vectors equal to  $(1, \dots, 1) \in \mathbb{R}^{A_i}$ . This leaves  $|A_{-i}| - 1$  linearly dependent columns. Eliminating them, genericity requires that for every  $i$ ,

$$|A_i| = |A_i| |S_i| \leq |A_{-i}| |S_{-i}| - (|A_{-i}| - 1) = |A_{-i}| \times (|S_{-i}| - 1) + 1.$$

Similarly, there are  $|A_i| - 1$  redundant row vectors when  $|S_{-i}| = 1$ . Since the intersection of finitely many generic sets is generic, DUD is generic if all these conditions hold.  $\square$

*Corollary 2.* By the Alternative Theorem,  $\Pr$  satisfies DUD $_B^B$  if and only if it satisfies PSI $_B^B$ , i.e., there exists a scheme  $\xi : I \times A \times S \rightarrow \mathbb{R}$  such that  $\xi_i(a, s) = 0$  if  $a \notin B$  and

$$\forall i \in I, a_i \in B_i, b_i \in A_i, \rho_i \in R_i, \quad 0 \leq \sum_{(a_{-i}, s)} \xi_i(a, s) (\Pr(s|a_{-i}, b_i, \rho_i) - \Pr(s|a)),$$

with a strict inequality whenever  $a_i \neq b_i$ . Replacing  $\xi_i(a, s) = \sigma(a)\zeta_i(a, s)$  for some (or equivalently any) correlated strategy  $\sigma$  with  $\text{supp } \sigma = B$ , this is equivalent to there being, for every profile of utility functions, an appropriate rescaling of  $\zeta$  that satisfies (\*).  $\square$

*Theorem 3.* For necessity, if  $\text{DUD}_B$  fails then there is a  $B$ -disobedient, undetectable deviation plan  $\alpha_i$ . Therefore,  $\alpha_i(b_i, \rho_i|a_i) > 0$  for some  $a_i \in B_i$ ,  $b_i \neq a_i$  and  $\rho_i \in R_i$ . Letting  $v_i(a_{-i}, b_i) < v_i(a)$  for every  $a_{-i}$ , clearly no correlated strategy with positive probability on  $a_i$  is approximately enforceable. Sufficiency follows by Lemmata B.2, B.3 and B.9.  $\square$

*Theorem 4.* Consider the family of linear programs below indexed by  $z \in [0, \infty)^I$ .

$$\begin{aligned} \max_{\varepsilon \geq 0, \xi} \quad & - \sum_{(i, a_i)} \varepsilon_i(a_i) \quad \text{s.t.} \quad \forall(i, a, s), \quad -\sigma(a)z_i \leq \xi_i(a, s) \leq \sigma(a)z_i, \\ & \forall(i, a_i, b_i, \rho_i), \quad \sum_{a_{-i}} \sigma(a) \Delta v_i(a, b_i) - \sum_{a_{-i}} \xi_i(a) \cdot \Delta \Pr(a, b_i, \rho_i) \leq \varepsilon_i(a_i), \end{aligned}$$

where  $\Delta v_i(a, b_i) := v_i(a_{-i}, b_i) - v_i(a)$  and  $\Delta \Pr(a, b_i, \rho_i) := \Pr(a_{-i}, b_i, \rho_i) - \Pr(a)$ . Given  $z \geq 0$ , the primal problem above looks for a scheme  $\xi$  adapted to  $\sigma$  (i.e., such that  $\xi_i(a, s) = 0$  whenever  $\sigma(a) = 0$ ) that minimizes the burden  $\varepsilon_i(a_i)$  of relaxing incentive constraints. By construction,  $\sigma$  is enforceable with transfers bounded by  $z$  if and only if there is a feasible  $\xi$  with  $\varepsilon_i(a_i) = 0$  for all  $(i, a_i)$ , i.e., the value of the problem is zero. Since  $\sigma$  is assumed enforceable, such  $z$  exists. The dual of this problem is:

$$\begin{aligned} \min_{\alpha, \beta \geq 0} \quad & \sum_{(i, a)} \sigma(a) [z_i \sum_{s \in S} \sigma(a) (\beta_i^+(a, s) + \beta_i^-(a, s)) - \Delta v_i(a, \alpha_i)] \quad \text{s.t.} \\ & \forall(i, a_i), \quad \sum_{(b_i, \rho_i)} \alpha_i(b_i, \rho_i|a_i) \leq 1, \\ & \forall i \in I, a \in \text{supp } \sigma, s \in S, \quad \Delta \Pr(s|a, \alpha_i) = \beta_i^+(a, s) - \beta_i^-(a, s). \end{aligned}$$

Since  $\beta_i^\pm(a, s) \geq 0$ , it follows easily that  $\beta_i^+(a, s) = \max\{\Delta \Pr(s|a, \alpha_i), 0\}$  and  $\beta_i^-(a, s) = \min\{\Delta \Pr(s|a, \alpha_i), 0\}$ . Hence,  $\beta_i^+(a, s) + \beta_i^-(a, s) = |\Delta \Pr(s|a, \alpha_i)|$ . Since  $\|\Delta \Pr(a, \alpha_i)\| = \sum_s |\Delta \Pr(s|a, \alpha_i)|$ , the dual is now equivalent to

$$V_\sigma(z) = \max_{\alpha \geq 0} \sum_{(i, a)} \sigma(a) (\Delta v_i(a, \alpha_i) - z \|\Delta \Pr(a, \alpha_i)\|) \quad \text{s.t.} \quad \forall(i, a_i), \quad \sum_{(b_i, \rho_i)} \alpha_i(b_i, \rho_i|a_i) \leq 1.$$

Adding mass to  $\alpha_i(a_i, \tau_i|a_i)$  if necessary, without loss  $\alpha_i$  is a deviation plan, proving (i).

To prove (ii), the first sentence is obvious. The second follows by Corollary 3: if  $\sigma$  is not enforceable then a  $\sigma$ -profitable,  $\text{supp } \sigma$ -undetectable plan  $\alpha_i$  exists, so  $V_\sigma(z) > 0$  for all  $z$ .

For (iii), if  $\sigma$  is not enforceable then there is a  $\sigma$ -profitable,  $\text{supp } \sigma$ -undetectable deviation plan  $\alpha_i^*$ . Approaching  $\alpha_i^*$  from  $\mathcal{F}_i$  (e.g., with mixtures of  $\alpha_i^*$  and a fixed plan in  $\mathcal{F}_i$ ), the

denominator defining  $\bar{z}_i$  tends to zero whilst the numerator tends to a positive amount, so  $\bar{z}_i$  is unbounded. Conversely, suppose  $\sigma$  is enforceable. If the sup defining  $\bar{z}_i$  is attained, we are done. If not, it is approximated by a sequence of supp  $\sigma$ -detectable deviation plans that converge to a supp  $\sigma$ -undetectable one. Since  $\sigma$  is enforceable, the limit is unprofitable. Let

$$F_i^\sigma(\delta) := \min_{\lambda_i \geq 0} \sum_{a \in A} \sigma(a) \|\Delta \Pr(a, \lambda_i)\| \quad \text{s.t.} \quad \Delta v_i(\sigma, \lambda_i) \geq \delta.$$

Since every  $\sigma$ -profitable deviation plan is detectable by [Corollary 3](#), it follows that  $F_i^\sigma(\delta) > 0$  for all  $\delta > 0$ , and  $\bar{z}_i = (\lim_{\delta \downarrow 0} F_i^\sigma(\delta)/\delta)^{-1}$ . Hence, it suffices to show  $\lim_{\delta \downarrow 0} F_i^\sigma(\delta)/\delta > 0$ . To this end, by adding variables like  $\beta$  above, the dual problem for  $F_i^\sigma$  is equivalent to:

$$\begin{aligned} F_i^\sigma(\delta) = \max_{\varepsilon \geq 0, x_i} \quad & \varepsilon \delta \quad \text{s.t.} \quad \forall(a, s), \quad -1 \leq x_i(a, s) \leq 1, \\ \forall(a_i, b_i, \rho_i), \quad & \sum_{a-i} \sigma(a) (\varepsilon \Delta v_i(a, b_i) - x_i(a) \cdot \Delta \Pr(a, b_i, \rho_i)) \leq 0. \end{aligned}$$

Since  $\sigma$  is enforceable, there is a feasible solution to this dual  $(\varepsilon, x_i)$  with  $\varepsilon > 0$ . Hence,  $F_i^\sigma(\delta) \geq \varepsilon \delta$  for all  $\delta > 0$ , therefore  $\lim_{\delta \downarrow 0} F_i^\sigma(\delta)/\delta > 0$ , as claimed.

To prove (iv), suppose that  $\bar{z}_i < \infty$  for all  $i$ . We claim  $V_\sigma(\bar{z}) = 0$ . Indeed, given  $\alpha_i^* \in \mathcal{F}_i$  for all  $i$ , substituting the definition of  $\bar{z}_i$  into the objective of the minimization in (i),

$$\sum_{i \in I} \Delta v_i(\sigma, \alpha_i^*) - \sum_{(i,a)} \sigma(a) \sup_{\alpha_i \in \mathcal{F}_i} \left\{ \frac{\max\{\Delta v_i(\sigma, \alpha_i), 0\}}{\sum_a \sigma(a) \|\Delta \Pr(a, \alpha_i)\|} \right\} \|\Delta \Pr(a, \alpha_i^*)\| \leq 0.$$

If  $\alpha_i^* \notin \mathcal{F}_i$  then, since  $\sigma$  is enforceable, every supp  $\sigma$ -undetectable deviation plan is unprofitable, so again the objective is non-positive, hence  $V_\sigma(\bar{z}) = 0$ . Clearly,  $V_\sigma$  decreases with  $z$ , so it remains to show that  $V_\sigma(\bar{z}) > 0$  if  $z_i < \bar{z}_i$  for some  $i$ . But by definition of  $\bar{z}$ , there is a deviation plan  $\alpha_i^*$  with  $\Delta v_i(\sigma, \alpha_i^*)/\sum_a \sigma(a) \|\Delta \Pr(a, \alpha_i^*)\| > z_i$ , so  $V_\sigma(z) > 0$ .  $\square$

**Theorem 5.** For sufficiency, suppose that  $\sigma$  is approximately enforceable, so there is a sequence  $\{\sigma^m\}$  such that  $\sigma^m$  is enforceable for every  $m$  and  $\sigma^m \rightarrow \sigma$ . Without loss, assume that  $\text{supp } \sigma^m \supset \text{supp } \sigma$  for all  $m$ . If  $\sigma^m = \sigma$  for all large  $m$  then  $\sigma$  is enforceable and the condition of [Theorem 5](#) is fulfilled with  $\mu = \sigma$ , so suppose not. If there exists  $m$  and  $m'$  such that  $\sigma^m = p\sigma^{m'} + (1-p)\sigma$  then incentive compatibility with respect to  $m$  yields that  $\sum_{a-i} \sigma^m(a) \Delta v_i(a, \alpha_i) \leq \sum_{a-i} \sigma^m(a) \zeta_i^m(a) \cdot \Delta \Pr(a, \alpha_i) \leq \sum_{a-i} \sigma^m(a) \bar{z} \|\Delta \Pr(a, \alpha_i)\|$  for every  $\alpha_i$ , where  $\bar{z} = \max_{(i,a,s)} |\zeta_i^m(a, s)|$ . For large  $m'$ ,  $\sigma^{m'}$  is sufficiently close to  $\sigma$  that if  $\alpha_i$  is  $\sigma$ -profitable then  $\sum_{a-i} \sigma^{m'}(a) \Delta v_i(a, \alpha_i) > 0$ , so  $\alpha_i$  is detectable.

If there does not exist  $m$  and  $m_1$  such that  $\sigma^m = p\sigma^{m_1} + (1-p)\sigma$  then there exists  $\sigma^{m_2}$  such that its distance from  $\sigma$  is less than the positive minimum distance between  $\sigma$  and the affine hull of  $\{\sigma^m, \sigma^{m_1}\}$ . Therefore, the lines generated by  $\sigma^m$  and  $\sigma^{m_1}$  and  $\sigma^{m_1}$  and  $\sigma^{m_2}$  are not



collinear. Proceeding inductively, pick  $C = \{\sigma^{m_1}, \dots, \sigma^{m_{|A|}}\}$  such that its affine space is full-dimensional in  $\Delta(A)$ . Since we are assuming that  $\sigma$  is not enforceable, it lies outside  $\text{conv } C$ . Let  $\hat{\sigma} = \sum_k \sigma^{m_k} / |A|$  and  $B_\varepsilon(\hat{\sigma})$  be the open  $\varepsilon$ -ball around  $\hat{\sigma}$  for some  $\varepsilon > 0$ . By construction,  $B_\varepsilon(\hat{\sigma}) \subset \text{conv } C$  for  $\varepsilon > 0$  sufficiently small, so there exists  $\hat{\sigma}' \in B_\varepsilon(\hat{\sigma})$  such that  $p\hat{\sigma} + (1-p)\sigma = \hat{\sigma}'$  for some  $p$  such that  $0 < p < 1$ . Now we can apply the argument from the previous paragraph, so the condition of [Theorem 5](#) holds.

For necessity, if  $\sigma$  is not approximately enforceable then  $1 \geq V_\sigma(z) \geq C > 0$  for every  $z$ , where  $V_\sigma$  is defined in [Lemma B.2](#). Let  $(\lambda^z, \mu^z)$  solve  $V_\sigma(z)$  for every  $z$ . Given  $\mu \in \Delta(A)$ ,

$$C \leq V_\sigma(z) \leq 1 + \sum_{(i,a)} \Delta v_i(\mu, \lambda_i^z) - z \sum_{(i,a)} \mu(a) \|\Delta \text{Pr}(a, \lambda_i^z)\|.$$

If the condition of [Theorem 5](#) holds then  $\sum_{(i,a)} \Delta v_i(\mu, \lambda_i^z) < \bar{z} \sum_{(i,a)} \mu(a) \|\Delta \text{Pr}(a, \lambda_i^z)\|$  and  $\sum_{(i,a)} \mu(a) \|\Delta \text{Pr}(a, \lambda_i^z)\| > 0$ , since there must exist  $i$  such that  $\lambda_i^\sigma$  is  $\sigma$ -profitable. Hence,  $C \leq 1 + (\bar{z} - z) \sum_{(i,a)} \mu(a) \|\Delta \text{Pr}(a, \lambda_i^z)\|$ , i.e.,  $z - \bar{z} \leq (1 - c) / \sum_{(i,a)} \mu(a) \|\Delta \text{Pr}(a, \lambda_i^z)\|$ . This inequality must hold for every  $z$ , therefore  $\sum_{(i,a)} \mu(a) \|\Delta \text{Pr}(a, \lambda_i^z)\| \rightarrow 0$  as  $z \rightarrow \infty$ . But this contradicts [Lemma B.10](#), since  $\sum_i \Delta v_i(\sigma, \lambda_i) \geq C$ , completing the proof.  $\square$

## B Lemmata

**Lemma B.1.** *A monitoring technology satisfies DUD if*

$$\forall(i, a_i, s_i), \quad \text{Pr}(a_i, s_i) \notin \text{cone}\{\text{Pr}(b_i, t_i) : (b_i, t_i) \neq (a_i, s_i)\},$$

where cone stands for the set of positive linear combinations of  $\{\text{Pr}(b_i, t_i) : (b_i, t_i) \neq (a_i, s_i)\}$ .

*Proof.* If DUD fails then there exists  $\alpha_i$  such that  $\alpha_i(b_i, \rho_i|a_i) > 0$  for some  $a_i \neq b_i$  and

$$\begin{aligned} \forall(a, s), \quad \text{Pr}(s|a) &= \sum_{(b_i, \rho_i)} \sum_{t_i \in \rho_i^{-1}(s_i)} \alpha_i(b_i, \rho_i|a_i) \text{Pr}(s_{-i}, t_i|a_{-i}, b_i) \\ &= \sum_{(b_i, t_i)} \sum_{\{\rho_i: \rho_i(t_i)=s_i\}} \alpha_i(b_i, \rho_i|a_i) \text{Pr}(s_{-i}, t_i|a_{-i}, b_i). \end{aligned}$$

Write  $\lambda_i(a_i, s_i, b_i, t_i) := \sum_{\{\rho_i: \rho_i(t_i)=s_i\}} \alpha_i(b_i, \rho_i|a_i)$ . By construction,  $\lambda_i(a_i, s_i, b_i, t_i) \geq 0$  is strictly positive for some  $a_i \neq b_i$  and satisfies

$$\forall(i, a, s), \quad \text{Pr}(s|a) = \sum_{(b_i, t_i)} \lambda_i(a_i, s_i, b_i, t_i) \text{Pr}(s_{-i}, t_i|a_{-i}, b_i).$$

Without loss,  $\lambda_i(a_i, s_i, a_i, s_i) = 0$  for some  $(a_i, s_i)$ . Indeed, if  $\lambda_i(a_i, s_i, a_i, s_i) = 1$  for all  $(a_i, s_i)$ , then the equation above is violated because  $\alpha_i$  is disobedient by hypothesis and

probabilities are non-negative. If  $\lambda_i(a_i, s_i, a_i, s_i) \neq 1$  then subtract  $\lambda_i(a_i, s_i, a_i, s_i) \Pr(s|a)$  from both sides of the equation and divide by  $1 - \lambda_i(a_i, s_i, a_i, s_i)$ . Therefore,  $\Pr(a_i, s_i) \in \text{cone}\{\Pr(b_i, t_i) : (b_i, t_i) \neq (a_i, s_i)\}$  for some  $(a_i, s_i)$ .  $\square$

**Lemma B.2.** *Consider the following linear program.*

$$\begin{aligned} V_\sigma(z) &:= \min_{\mu \geq 0, p, \xi} p \quad \text{s.t.} \quad \sum_{a \in A} \mu(a) = p, \\ \forall(i, a, s), \quad & -(\mu(a) + (1-p)\sigma(a))z \leq \xi_i(a, s) \leq (\mu(a) + (1-p)\sigma(a))z, \\ \forall(i, a_i, b_i, \rho_i), \quad & \sum_{a-i} (\mu(a) + (1-p)\sigma(a)) \Delta v_i(a, b_i) \leq \sum_{a-i} \xi_i(a) \cdot \Delta \Pr(a, b_i, \rho_i). \end{aligned}$$

The correlated strategy  $\sigma$  is approximately enforceable if and only if  $V_\sigma(z) \rightarrow 0$  as  $z \rightarrow \infty$ . The dual of the above linear program is given by the following problem:

$$\begin{aligned} V_\sigma(z) &= \max_{\lambda \geq 0, \kappa} \sum_{i \in I} \Delta v_i(\sigma, \lambda_i) - z \sum_{(i,a)} \sigma(a) \|\Delta \Pr(a, \lambda_i)\| \quad \text{s.t.} \\ \forall a \in A, \quad & \kappa \leq \sum_{i \in I} \Delta v_i(a, \lambda_i) - z \sum_{i \in I} \|\Delta \Pr(a, \lambda_i)\|, \\ & \sum_{i \in I} \Delta v_i(\sigma, \lambda_i) - z \sum_{(i,a)} \sigma(a) \|\Delta \Pr(a, \lambda_i)\| = 1 + \kappa. \end{aligned}$$

*Proof.* The first family of primal constraints require  $\xi$  to be adapted to  $\mu + (1-p)\sigma$ , so for any  $z$ ,  $(\mu, p, \xi)$  solves the primal if and only if  $\mu + (1-p)\sigma$  is exactly enforceable with  $\xi$ . (Since correlated equilibrium exists, the primal constraint set is clearly nonempty, and for finite  $z$  it is also clearly bounded). The first statement now follows. The second statement follows by a lengthy but standard manipulation of the primal to obtain the above dual.  $\square$

**Lemma B.3.** *Consider the following family of linear programs indexed by  $\varepsilon > 0$  and  $z \geq 0$ .*

$$\begin{aligned} F_\sigma^\varepsilon(z) &:= \max_{\lambda \geq 0} \min_{\mu \in \Delta(A)} \sum_{i \in I} \Delta v_i(\mu, \lambda_i) - z \sum_{(i,a)} \mu(a) \|\Delta \Pr(a, \lambda_i)\| \quad \text{s.t.} \\ & \sum_{i \in I} \Delta v_i(\sigma, \lambda_i) - z \sum_{(i,a)} \sigma(a) \|\Delta \Pr(a, \lambda_i)\| \geq \varepsilon. \end{aligned}$$

$F_\sigma^\varepsilon(z) \rightarrow -\infty$  as  $z \rightarrow \infty$  for some  $\varepsilon > 0$  if and only if  $\sigma$  is approximately enforceable.

*Proof.* The dual of the problem defining  $F_\sigma^\varepsilon(z)$  is

$$\begin{aligned} F_\sigma^\varepsilon(z) &= \min_{\delta, \mu \geq 0, x} -\delta \varepsilon \quad \text{s.t.} \quad \sum_{a \in A} \mu(a) = 1, \\ \forall(i, a, s), \quad & -(\mu(a) + \delta \sigma(a))z \leq x_i(a, s) \leq (\mu(a) + \delta \sigma(a))z, \\ \forall(i, a_i, b_i, \rho_i), \quad & \sum_{a-i} (\mu(a) + \delta \sigma(a)) \Delta v_i(a, b_i) \leq \sum_{a-i} x_i(a) \cdot \Delta \Pr(a, b_i, \rho_i). \end{aligned}$$

Since clearly  $\varepsilon > 0$  does not affect the dual feasible set, if  $F_\sigma^\varepsilon(z) \rightarrow -\infty$  for some  $\varepsilon > 0$  then there exists  $z \geq 0$  such that  $\delta > 0$  is feasible, and  $\delta \rightarrow \infty$  as  $z \rightarrow \infty$ . Therefore,  $F_\sigma^\varepsilon(z) \rightarrow -\infty$  for every  $\varepsilon > 0$ . If  $V_\sigma(z) = 0$  for some  $z$  we are done by monotonicity of  $V_\sigma$ . Otherwise, suppose that  $V_\sigma(z) > 0$  for all  $z > 0$ . Let  $(\lambda, \kappa)$  be an optimal dual solution for  $V_\sigma(z)$  in Lemma B.2. By optimality,  $\kappa = \min_{\mu \in \Delta(A)} \sum_i \Delta v_i(\mu, \lambda_i) - z \sum_{(i,a)} \mu(a) \|\Delta \Pr(a, \lambda_i)\|$ . Therefore, by the second dual constraint in  $V_\sigma(z)$  of Lemma B.2,

$$V_\sigma(z) = 1 + \kappa = 1 + F_\sigma^{V_\sigma(z)}(z) = 1 - \delta V_\sigma(z),$$

where  $\delta$  is an optimal solution to the dual with  $\varepsilon = V_\sigma(z)$ . Rearranging,  $V_\sigma(z) = 1/(1 + \delta)$ . Finally,  $F_\sigma^\varepsilon(z) \rightarrow -\infty$  as  $z \rightarrow \infty$  if and only if  $\delta \rightarrow \infty$ , if and only if  $V_\sigma(z) \rightarrow 0$ .  $\square$

**Lemma B.4.** Fix any  $\varepsilon > 0$ . If  $\Pr$  satisfies  $DUD_B$ , where  $B = \text{supp } \sigma$ , then for every  $C \leq 0$  there exists  $z \geq 0$  such that  $G_\sigma(z) \leq C$ , where  $\Delta v_i(a_i)^* := \max_{(a_i, b_i)} \{\Delta v_i(a, b_i)\}$ ,

$$\Delta v_i(a_i, \lambda_i)^* := \Delta v_i(a_i)^* \sum_{(a_i, b_i \neq a_i, \rho_i)} \lambda_i(a_i, b_i, \rho_i), \text{ and}$$

$$G_\sigma(z) := \max_{\lambda \geq 0} \sum_{(i,a)} \|\Delta v_i(a_i, \lambda_i)\| - z \sum_{(i,a)} \|\Delta \Pr(a, \lambda_i)\| \quad \text{s.t.}$$

$$\forall i \in I, a_i \notin B_i, \lambda_i(a_i) = 0, \quad \text{and} \quad \sum_{i \in I} \Delta v_i(\sigma, \lambda_i) - z \sum_{(i,a)} \sigma(a) \|\Delta \Pr(a, \lambda_i)\| \geq \varepsilon.$$

*Proof.* The dual of this problem is given by

$$\begin{aligned} G_\sigma(z) = \min_{\delta \geq 0, x} & -\delta \varepsilon \quad \text{s.t.} \\ \forall (i, a, s), & -(1 + \delta \sigma(a))z \leq x_i(a, s) \leq (1 + \delta \sigma(a))z, \\ \forall (i, a_i \in B_i, b_i, \rho_i), & \sum_{a_{-i}} \delta \sigma(a) \Delta v_i(a, b_i) + \mathbf{1}_{\{a_i \neq b_i\}} \Delta v_i(a_i)^* \leq \sum_{a_{-i}} x_i(a) \cdot \Delta \Pr(a, b_i, \rho_i), \end{aligned}$$

where  $\mathbf{1}_{\{b_i \neq a_i\}} = 1$  if  $b_i \neq a_i$  and 0 otherwise. This problem looks almost exactly like the dual for  $F_\sigma^\varepsilon(z)$  except that the incentive constraints are only indexed by  $a_i \in B_i$ . Now,  $DUD_B$  is equivalent to  $PSI_B$ , i.e., there is an incentive scheme  $x : I \times A \times S \rightarrow \mathbb{R}$  such that

$$\forall (i, a_i, b_i, \rho_i), \quad 0 \leq \sum_{a_{-i}} x_i(a) \cdot \Delta \Pr(a, b_i, \rho_i)$$

with a strict inequality whenever  $a_i \in B_i$  and  $a_i \neq b_i$ . Hence, by scaling  $x$  appropriately, there is a feasible dual solution with  $\delta > 0$ , so  $G_\sigma(z) < 0$ . Moreover, for any  $\delta > 0$ , clearly an  $x$  exists with  $\sum_{a_{-i}} \delta \sigma(a) \Delta v_i(a, b_i) + \mathbf{1}_{\{b_i \neq a_i\}} \Delta v_i(a_i)^* \leq \sum_{a_{-i}} x_i(a) \cdot \Delta \Pr(a, b_i, \rho_i)$  on all  $(i, a_i \in B_i, b_i, \rho_i)$  by  $PSI_B$ , so there exists  $z$  to make such  $\delta$  feasible. In particular,  $\delta \geq C/\varepsilon$  is feasible for some  $z$ , as required.  $\square$

**Lemma B.5.** If  $\Pr$  satisfies  $DUD_B$ , then there exists a finite  $z \geq 0$  such that

$$\forall i \in I, a_i \in B_i, \lambda_i \geq 0, \quad \sum_{a_{-i}} \Delta v_i(a_i, \lambda_i)^* - z \|\Delta \Pr(a, \lambda_i)\| \leq 0.$$

*Proof.* Given  $i$ ,  $a_i \in B_i$ , plug  $\sigma(a) = 1/|A_{-i}|$  for all  $a_{-i}$  in the proof of [Theorem 4 \(iii\)](#).  $\square$

Call  $\lambda$  *extremely detectable* if  $\lambda_i(a_i)$  cannot be written as a positive linear combination involving undetectable deviations (possibly mixed) for every  $(i, a_i)$ . Let  $\mathcal{E}$  denote the set of all such extremely detectable  $\lambda$ .

**Lemma B.6.** *The set  $\mathcal{D}^e = \{\alpha \in \mathcal{E} : \forall(i, a_i), \sum_{(b_i, \rho_i)} \alpha_i(a_i, b_i, \rho_i) = 1\}$  is compact.*

*Proof.*  $\mathcal{D}^e$  is clearly a bounded subset of Euclidean space, so it remains to show that it is closed. Consider a sequence  $\{\alpha^m\} \subset \mathcal{D}^e$  such that  $\alpha^m \rightarrow \alpha^*$ . For any  $\alpha \in \mathcal{D}$ , let

$$p^*(\alpha) := \max_{0 \leq p \leq 1, \alpha^i \in \mathcal{D}} \{p : \alpha^0 \text{ is undetectable, } p\alpha^0 + (1-p)\alpha^1 = \alpha\}.$$

This is a well-defined linear program with a compact constraint set and finite values, so  $p^*$  is continuous in  $\alpha$ . By assumption,  $p^*(\alpha^m) = 0$  for all  $m$ , so  $p^*(\alpha^*) = 0$ , hence  $\alpha^* \in \mathcal{D}^e$ .  $\square$

**Lemma B.7.** *Let  $\mathcal{D}^e$  be the set of extremely detectable deviation plans.*

$$\gamma := \min_{\alpha^e \in \mathcal{D}^e} \sum_{(i, a)} \|\Delta \Pr(a, \alpha_i^e)\| > 0.$$

*Proof.* If  $\mathcal{D}^e = \emptyset$  then  $\gamma = +\infty$ . If not,  $\mathcal{D}^e$  is compact by [Lemma B.6](#), so there is no sequence  $\{\alpha_i^{e, m}\} \subset \mathcal{D}^e$  with  $\|\Delta \Pr(a, \alpha_i^{e, m})\| \rightarrow 0$  for all  $(i, a)$  as  $m \rightarrow \infty$ , hence  $\gamma > 0$ .  $\square$

**Lemma B.8.** *Let  $\mathcal{D}_i^e = \text{proj}_i \mathcal{D}^e$ . There exists a finite  $z \geq 0$  such that*

$$\forall i \in I, a_i \notin B_i, \alpha_i^e \in \mathcal{D}_i^e, \quad \sum_{a_{-i}} \Delta v_i(a_i, \alpha_i^e)^* - z \|\Delta \Pr(a, \alpha_i^e)\| \leq 0.$$

*Proof.* Let  $\|\Delta v\| = \max_{(i, a, b_i)} |\Delta v_i(a, b_i)|$ . If  $z \geq \|\Delta v\|/\gamma$ , with  $\gamma$  as in [Lemma B.7](#), then

$$\forall(i, a_i), \quad \sum_{a_{-i}} \Delta v_i(a_i, \alpha_i^e)^* - z \|\Delta \Pr(a, \alpha_i^e)\| \leq \|\Delta v\| - z \sum_{a_{-i}} \|\Delta \Pr(a, \alpha_i^e)\| \leq \|\Delta v\| - \frac{\|\Delta v\|}{\gamma} \gamma.$$

The right-hand side clearly equals zero, which establishes the claim.  $\square$

**Lemma B.9.** *Fix any  $\varepsilon > 0$ . If  $\Pr$  satisfies  $DUD_B$  then for every  $C \leq 0$  there exists  $z \geq 0$  such that for every  $\lambda \geq 0$  with*

$$\sum_{i \in I} \Delta v_i(\sigma, \lambda_i) - z \sum_{(i, a)} \sigma(a) \|\Delta \Pr(a, \lambda_i)\| \geq \varepsilon,$$

*there exists  $\mu \in \Delta(A)$  such that*

$$W(\mu, \lambda) := \sum_{i \in I} \Delta v_i(\mu, \lambda_i) - z \sum_{(i, a)} \mu(a) \|\Delta \Pr(a, \lambda_i)\| \leq C.$$

*Proof.* Rewrite  $W(\mu, \lambda)$  by splitting it into three parts,  $W_d(\mu, \lambda)$ ,  $W_e(\mu, \lambda)$  and  $W_u(\mu, \lambda)$ :

$$\begin{aligned} W_d(\mu, \lambda) &= \sum_{i \in I} \sum_{a_i \in B_i} \sum_{a_{-i}} \mu(a) (\Delta v_i(a, \lambda_i) - z \|\Delta \Pr(a, \lambda_i)\|) \\ W_e(\mu, \lambda) &= \sum_{i \in I} \sum_{a_i \notin B_i} \sum_{a_{-i}} \mu(a) (\Delta v_i(a, \lambda_i^e) - z \|\Delta \Pr(a, \lambda_i^e)\|), \\ W_u(\mu, \lambda) &= \sum_{i \in I} \sum_{a_i \notin B_i} \sum_{a_{-i}} \mu(a) (\Delta v_i(a, \lambda_i^u) - z \|\Delta \Pr(a, \lambda_i^u)\|), \end{aligned}$$

and  $\lambda = \lambda^e + \lambda^u$  with  $\lambda^e$  extremely detectable,  $\lambda^u$  undetectable. Since  $\lambda^u$  is undetectable,

$$W_u(\mu, \lambda) = \sum_{i \in I} \sum_{a_i \notin B_i} \sum_{a_{-i}} \mu(a) \Delta v_i(a, \lambda_i^u)$$

Let  $\mu^0(a) = 1/|A|$  for every  $a$ . By [Lemma B.4](#), there exists  $z$  with  $W_d(\mu^0, \lambda) \leq C$  for every  $\lambda$ , and by [Lemma B.8](#) there exists  $z$  with  $W_e(\mu^0, \lambda) \leq 0$  for every  $\lambda$ . Therefore, if  $W_u(\mu^0, \lambda) \leq 0$  we are done. Otherwise, for every  $i$  and  $a_i, b_i \in A_i$ , let  $\mu_i^0(a_i) = 1/|A_i|$  and

$$\mu_i^1(b_i) := \sum_{(a_i, \rho_i)} \frac{\lambda_i^u(a_i, b_i, \rho_i)}{\sum_{(b'_i, \rho'_i)} \lambda_i^u(a_i, b'_i, \rho'_i)} \mu_i^0(a_i)$$

Iterate this rule to obtain a sequence  $\{\mu_i^m\}$  with limit  $\mu_i^\infty \in \Delta(A_i)$ . By construction,  $\mu_i^\infty$  is a  $\lambda_i^u$ -stationary distribution ([Nau and McCardle, 1990](#); [Myerson, 1997](#)). Therefore, given any  $a_{-i}$ , the deviation gains for every agent equal zero, i.e.,

$$\sum_{(a_i, b_i, \rho_i)} \mu_i^\infty(a_i) \lambda_i^u(a_i, b_i, \rho_i) (v_i(a_{-i}, b_i) - v_i(a)) = 0.$$

Let  $\mu^m(a) := \prod_i \mu_i^m(a_i)$  for all  $m$ . By construction,  $W_u(\mu^\infty, \lambda^u) = 0$ . We will show that  $W_d(\mu^\infty, \lambda) \leq C$  and  $W_e(\mu^\infty, \lambda) \leq 0$ . To see this, notice firstly that, since  $\lambda_i^u$  is undetectable, for any other agent  $j \neq i$ , any  $\lambda_j \geq 0$  and every action profile  $a \in A$ ,

$$\|\Delta \Pr(a, \lambda_j)\| = \|\Delta \Pr(a, \lambda_i^u, \lambda_j)\| \leq \|\Delta \Pr(a, \hat{\lambda}_i^u, \lambda_j)\|,$$

where  $\hat{\lambda}_i^u(a_i, b_i, \tau_i) = \sum_{\rho_i} \lambda_i^u(a_i, b_i, \rho_i)$  and  $\hat{\lambda}_i^u(a_i, b_i, \rho_i) = 0$  for all  $\rho_i \neq \tau_i$ ,

$$\Delta \Pr(a, \lambda_i^u, \lambda_j) = \sum_{(b_j, \rho_j)} \lambda_j(a_j, b_j, \rho_j) \sum_{(b_i, \rho_i)} \lambda_i^u(a_i, b_i, \rho_i) (\Pr(a, b_i, \rho_i, b_j, \rho_j) - \Pr(a, b_i, \rho_i)),$$

and  $\Pr(s|a, b_i, \rho_i, b_j, \rho_j) = \sum_{t_j \in \rho_j^{-1}(s_j)} \Pr(s_{-j}, t_j | a_{-j}, b_j, b_i, \rho_i)$ . Secondly, notice that

$$\begin{aligned} \forall i \in I, a_i \in B_i, \quad & \sum_{a_{-i}} \mu^m(a) (\Delta v_i(a, \lambda_i) - z \|\Delta \Pr(a, \lambda_i)\|) \leq \\ & \mu_i^m(a_i) \sum_{a_{-i}} \mu_{-i}^m(a_{-i}) (\Delta v_i(a_i, \lambda_i)^* - z \|\Delta \Pr(a, \lambda_i)\|) \leq \\ & \mu_i^m(a_i) \sum_{a_{-i}} \mu_{-i}^0(a_{-i}) (\Delta v_i(a_i, \lambda_i)^* - z \|\Delta \Pr(a, \lambda_i)\|) \leq \\ & \sum_{a_{-i}} \mu^0(a) (\Delta v_i(a_i, \lambda_i)^* - z \|\Delta \Pr(a, \lambda_i)\|). \end{aligned}$$

Indeed, the first inequality is obvious. The second one follows by repeated application of the previously derived inequality  $\|\Delta \Pr(a, \lambda_i)\| \leq \|\Delta \Pr(a, \hat{\lambda}_j^u, \lambda_i)\|$  for each agent  $j \neq i$  separately  $m$  times. The third inequality follows because (i)  $\mu_i^m(a_i) \geq \mu_i^0(a_i)$  for all  $m$  and  $a_i \in B_i$ , since  $B_i$  is a  $\hat{\lambda}_i^u$ -absorbing set, and (ii)  $\sum_{a_{-i}} \Delta v_i(a_i, \lambda_i)^* - z \|\Delta \Pr(a, \lambda_i)\| \leq 0$  for every  $(i, a_i)$  by [Lemma B.5](#). Therefore,  $W_d(\mu^\infty, \lambda) \leq W_d(\mu^m, \lambda) \leq W_d(\mu^0, \lambda) \leq C$ . Thirdly,

$$\begin{aligned} \forall i \in I, a_i \notin B_i, \quad & \sum_{a_{-i}} \mu_{-i}^m(a_{-i}) (\Delta v_i(a, \lambda_i^e) - z \|\Delta \Pr(a, \lambda_i^e)\|) \leq \\ & \sum_{a_{-i}} \mu_{-i}^m(a_{-i}) (\Delta v_i(a_i, \lambda_i^e)^* - z \|\Delta \Pr(a, \lambda_i^e)\|) \leq \\ & \sum_{a_{-i}} \mu_{-i}^0(a_{-i}) (\Delta v_i(a_i, \lambda_i^e)^* - z \|\Delta \Pr(a, \lambda_i^e)\|) \leq 0. \end{aligned}$$

The first inequality is again obvious, the second inequality follows by repeated application of  $\|\Delta \Pr(a, \lambda_i)\| \leq \|\Delta \Pr(a, \hat{\lambda}_j^u, \lambda_i)\|$ , and the third one follows from [Lemma B.8](#). Hence,  $W_e(\mu^m, \lambda) \leq 0$  for every  $m$ , therefore  $W_e(\mu^\infty, \lambda) \leq 0$ . This completes the proof. (This proof extends [Nau and McCardle \(1990\)](#) and [Myerson \(1997\)](#) by including transfers.)  $\square$

**Lemma B.10.** *The conditions of [Theorem 5](#) imply that for every  $\varepsilon > 0$  there exists  $\delta > 0$  such that  $\sum_i \Delta v_i(\sigma, \lambda_i) \geq \varepsilon$  implies that  $\sum_{(i,a)} \mu(a) \|\Delta \Pr(a, \lambda_i)\| \geq \delta$  for some  $\mu \in \Delta(A)$  with  $\sum_i \Delta v_i(\mu, \lambda_i) \leq \bar{z} \sum_{(i,a)} \mu(a) \|\Delta \Pr(a, \lambda_i)\|$ .*

*Proof.* Otherwise, there exists  $\varepsilon > 0$  such that for every  $\delta > 0$  some  $\lambda^\delta$  exists with  $\sum_i \Delta v_i(\sigma, \lambda_i^\delta) \geq \varepsilon$  but  $\sum_{(i,a)} \mu(a) \|\Delta \Pr(a, \lambda_i)\| < \delta$  whenever  $\mu \in \Delta(A)$  satisfies the given inequality  $\sum_i \Delta v_i(\mu, \lambda_i) \leq \bar{z} \sum_{(i,a)} \mu(a) \|\Delta \Pr(a, \lambda_i)\|$ . If  $\lambda^\delta$  is bounded for every  $\delta$  then  $\{\lambda^\delta\}$  has a convergent subsequence with limit  $\lambda^0$ . But this  $\lambda^0$  violates the conditions of [Theorem 5](#), so assume that  $\{\lambda^\delta\}$  is unbounded. A deviation plan  $\alpha_i^r$  is called *relatively undetectable* if  $\sum_{(i,a)} \mu(a) \|\Delta \Pr(a, \lambda_i)\| = 0$  whenever  $\mu \in \Delta(A)$  satisfies  $\sum_i \Delta v_i(\mu, \lambda_i) \leq \bar{z} \sum_{(i,a)} \mu(a) \|\Delta \Pr(a, \lambda_i)\|$ . Call  $\mathcal{D}_i^r$  the set of relatively undetectable plans. A deviation plan  $\alpha_i^s$  is called *relatively detectable* if

$$\max_{(p, \alpha_i, \alpha_i^r)} \{p : p\alpha_i^r + (1-p)\alpha_i = \alpha_i^s, \alpha_i \in \mathcal{D}_i, \alpha_i^r \in \mathcal{D}_i^r, p \in [0, 1]\} = 0.$$

Let  $\mathcal{D}_i^s$  be the set of relatively detectable plans. By the same argument as for [Lemma B.6](#),  $\mathcal{D}_i^s$  is a compact set, therefore, by the same argument as for [Lemma B.7](#),

$$\gamma_i^s := \min_{\alpha_i^s \in \mathcal{D}_i^s} \max_{\mu \in \Delta(A)} \left\{ \sum_{(i,a)} \mu(a) \|\Delta \Pr(a, \alpha_i^s)\| : \sum_{i \in I} \Delta v_i(\mu, \lambda_i) \leq \bar{z} \sum_{(i,a)} \mu(a) \|\Delta \Pr(a, \lambda_i)\| \right\} > 0.$$

Without loss,  $\lambda_i^\delta = \lambda_i^{r,\delta} + \lambda_i^{s,\delta}$ , where  $\lambda_i^{r,\delta}$  is relatively undetectable and  $\lambda_i^{s,\delta}$  is relatively detectable. By assumption,  $\lambda_i^{r,\delta}$  is  $\sigma$ -unprofitable, so  $\sum_{(b_i, \rho_i)} \lambda_i^{s,\delta}(a_i, b_i, \rho_i)$  is bounded below

by  $\beta > 0$ , say. (Otherwise,  $\sum_i \Delta v_i(\sigma, \lambda_i^\delta) < \varepsilon$  for small  $\delta > 0$ .) But this implies that

$$\max_{\mu \in \Delta(A)} \sum_{(i,a)} \mu(a) \left\| \Delta \Pr(a, \lambda_i^\delta) \right\| = \max_{\mu \in \Delta(A)} \sum_{(i,a)} \mu(a) \left\| \Delta \Pr(a, \lambda_i^{s,\delta}) \right\| \geq \beta \gamma_i^s > 0.$$

But this contradicts our initial assumption, which establishes the result.  $\square$

## References

- ALCHIAN, A. AND H. DEMSETZ (1972): “Production, Information Costs, and Economic Organization,” *American Economic Review*, 62, 777–795. 1
- CHE, Y.-K. AND J. KIM (2006): “Robustly Collusion-Proof Implementation,” *Econometrica*, 74, 1063–1107. 19
- CHESLOW, D. (2005): “Students, Bartenders Cited in Liquor Sting,” *The Daily Northwestern*, <http://www.dailynorthwestern.com/>. 1
- COMPTE, O. (1998): “Communication in Repeated Games with Imperfect Private Monitoring,” *Econometrica*, 66, 597–626. 13
- CREMER, J. AND R. MCLEAN (1988): “Full Extraction of the Surplus in Bayesian and Dominant Strategy Auctions,” *Econometrica*, 56, 1247–1257. 19
- D’ASPREMONT, C. AND L.-A. GÉRARD-VARET (1998): “Linear Inequality Methods to Enforce Partnerships under Uncertainty: An Overview,” *Games and Economic Behavior*, 25, 311–336. 13
- FUDENBERG, D., D. LEVINE, AND E. MASKIN (1994): “The Folk Theorem with Imperfect Public Information,” *Econometrica*, 62, 997–1039. 2, 13
- HOLMSTRÖM, B. (1982): “Moral Hazard in Teams,” *Bell Journal of Economics*, 13, 324–340. 4
- KANDORI, M. (2003): “Randomization, Communication, and Efficiency in Repeated Games with Imperfect Public Monitoring,” *Econometrica*, 71, 345–353. 13
- KANDORI, M. AND H. MATSUSHIMA (1998): “Private Observation, Communication, and Collusion,” *Econometrica*, 66, 627–652. 13
- LEGROS, P. AND H. MATSUSHIMA (1991): “Efficiency in Partnerships,” *Journal of Economic Theory*, 55, 296–322. 13

- LEGROS, P. AND S. MATTHEWS (1993): “Efficient and Nearly Efficient Partnerships,” *Review of Economic Studies*, 60, 599–611. 13
- MYERSON, R. (1986): “Multistage Games with Communication,” *Econometrica*, 54, 323–358. 2
- (1997): “Dual Reduction and Elementary Games,” *Games and Economic Behavior*, 21, 183–202. 28, 29
- NAU, R. F. AND K. F. MCCARDLE (1990): “Coherent Behavior in Noncooperative Games,” *Journal of Economic Theory*, 50, 424–444. 28, 29
- RAHMAN, D. (2005a): “Optimum Contracts with Public and Private Monitoring,” Mimeo.
- (2005b): “Team Formation and Organization,” Ph.D. dissertation, UCLA. 0
- ROCKAFELLAR, R. T. (1970): *Convex Analysis*, Princeton, New Jersey: Princeton University Press. 11
- TOMALA, T. (2005): “On Subgame-Perfect Communication Equilibria in Repeated Games with Imperfect Monitoring,” Mimeo. 14
- TSA (2004): “Guidance on Screening Partnership Program,” *Transportation Security Administration*, [www.tsa.gov/assets/pdf/ SPP\\_OptOut\\_Guidance\\_6.21.04.pdf](http://www.tsa.gov/assets/pdf/SPP_OptOut_Guidance_6.21.04.pdf). 1