

# Strategic Ambiguity and Arms Proliferation

Sandeep Baliga  
Northwestern University

Tomas Sjöström  
Rutgers University

November 8, 2006

## Abstract

A big power is facing a small power that may have developed WMDs. The small power can create *strategic ambiguity* by not allowing arms inspections. We study the impact of strategic ambiguity on arms proliferation and the probability of conflict. Creating strategic ambiguity is a substitute for actually acquiring new weapons: ambiguity reduces the incentive for the small power to invest in a weapons program, which reduces the risk of arms proliferation. Therefore, strategic ambiguity tends to benefit the big power. On the other hand, strategic ambiguity may hurt the small power because it does not always protect it from an attack. Cheap-talk messages can be used to trigger inspections when they are most valuable to the big power. To preserve incentive compatibility, the “tough” messages which make inspections more likely must imply a greater risk of arms proliferation.

## 1 Introduction

For several years, the North Korean regime claimed it had nuclear weapons, but it did not allow outsiders to verify this claim. There was *strategic ambiguity* about their capabilities (Norris and Kristensen [14]). North Korea’s Vice Minister of Foreign Affairs told visiting American scientists that the ambiguity protected the regime against punitive actions: “If you go back to the United States and say that the North already has nuclear weapons, this may cause the U.S. to act against us” (Hecker [9]). Thus, a regime which has advanced weapons may be better off not revealing it. Strategic ambiguity

can also benefit a regime which *lacks* advanced weapons, since the ambiguity might deter opportunistic enemies from attacking. This may have been Saddam Hussein's (failed) policy after 1991. There is wide-spread belief that the U.S. attacked Iraq only after it concluded that Iraq did *not* have WMDs. Iranian newspaper editorials argue that "In the contemporary world, it is obvious that having access to advanced weapons shall cause deterrence and therefore security, and will neutralize the evil wishes of great powers to attack other nations" (*Jumhuri-ye Islami*, cited by Takeyh [19]). Perhaps in truth North Korea and Iran have no need for nuclear deterrence. But their own perceptions, even if inaccurate, will determine their behavior.

It is often argued that strategic ambiguity makes the world a more dangerous place and, conversely, that weapons inspections make the world safer. According to this line of thinking, incomplete information about an opponent's capabilities generates fear which causes arms races and wars. Simply talking to the opponent may not reduce fear, because talk is cheap. However, if the opponent reveals that he is unarmed then fear is reduced. This standard argument is embodied in Article 3 of the *Treaty on the Non-proliferation of Nuclear Weapons*, known as the NPT [22]. The NPT requires that nations submit to inspection and verification of nuclear facilities by the IAEA, which is meant to promote peace and trust. However, Israel, India and Pakistan have not signed the NPT, North Korea has withdrawn from it, and Iran is close to violating it. The exact quality and quantity of WMDs in these countries is unknown. For example, Pakistan and India possess short-range nuclear weapons but it is unclear whether they have intercontinental ballistic missiles or are developing them (Norris and Kristensen [15], Norris, Kristensen and Handler [13]). Why wouldn't a leader who desires peace sign the NPT, disarm, and then allow arms inspectors to verify it? The problem is obvious: a leader who doesn't trust his opponent is unwilling to disarm and then reveal that he is defenseless. The NPT cannot create trust where no trust exists to begin with. Therefore, a withdrawal from the NPT is not necessarily a signal that the leader has aggressive intentions, neither is it necessarily a signal that a nuclear arsenal exists. Strategic ambiguity can be a *substitute* for WMDs. Ambiguity can provide deterrence and thereby reduce the value of acquiring WMDs, which leads to less arms proliferation. The standard argument in favor of weapons inspections does not take this into account.

Sobel [18] gives an argument in favor of ambiguity. A well-armed country is less likely to be attacked, so if disclosure is possible then all well-armed

countries prefer to reveal their weapons. In equilibrium there is no strategic ambiguity, because any country which does not reveal its capabilities will be known to be weak and may be attacked. (See Grossman [8] and Milgrom [11] for similar “unraveling” arguments.) But if disclosure is impossible, then countries which are not well-armed cannot be singled out for opportunistic attacks. This argument suggests that small powers (whose capabilities are uncertain) are better off *ex ante* if disclosure is impossible, but big powers (who are known to be armed) prefer to eliminate all ambiguity so they can attack the weak. In practice, it may be difficult to prevent disclosure. For example, during the cold war each nuclear power advertised its capabilities, and it is unclear how they could have been prevented from doing so. However, in the asymmetric post cold war world it is not necessarily in the interest of a small power to reveal that it has WMDs, since doing so might cause a big power to act against them. Formally, we will argue that “hard” information about capabilities can be a negative signal about “soft” information about an actor’s type. In this case, the unraveling argument fails and strategic ambiguity can be sustained in equilibrium even though disclosure is possible.

Sobel’s [18] argument in favor of ambiguity relies on the assumption that a country’s strength is exogenously given. If this assumption is dropped then the argument becomes quite different. Suppose small powers can acquire advanced weapons at a cost. In this case, ambiguity requires that small powers arm themselves with a positive probability which is determined in equilibrium, together with a positive probability that the big power attacks. In an equilibrium with ambiguity, some attacks will be “mistakes”, i.e., they could have been avoided, had the small power’s true strength been publicly known. But the equilibrium probability that the small power arms itself is smaller than it would be with full disclosure. Thus, *ambiguity reduces the risk of arms proliferation*, which is good for the big power. Moreover, if cheap-talk is allowed, then those leaders who are more likely to make “mistakes” can send a message which increases the probability of an inspection, and thereby reduce the likelihood of a mistake.

In our “arms proliferation game” there are two players, A and B, who are the leaders of countries A and B. Country A is a big power that is known to possess advanced weapons. Country B is a small power which initially is unarmed, i.e., it lacks advanced weapons. But B can try to acquire advanced weapons by making an investment. If this succeeds, then B will become armed. Player A thinks there is a small probability that B is a “crazy” type who would share his advanced weapons with terrorists. Player B’s

true type is *soft* (unverifiable) information. In Baliga and Sjöström [1], we studied how soft private information can trigger arms races and wars. In the current model, we will focus on whether revelation of *hard* information can promote peace and trust. Specifically, we assume weapons inspectors can verify whether or not B is armed (but not whether or not he is “crazy”). Player A’s decision is whether or not to attack B. The optimal decision depends on A’s preferences (his type) and his beliefs. At most three equilibria can exist in this game. There is always an equilibrium with *full ambiguity* where B never allows arms inspections, and a *full disclosure* equilibrium where he always does. A more interesting possibility is a *communication equilibrium* involving informative cheap-talk.

With full disclosure, B’s fear that A may attack him impels B to invest. With full ambiguity, B has less incentive to invest, which tends to make A better off. However, if A is an *opportunistic type*, then he would like to attack if B is unarmed, in order to obtain control of some resource or simply achieve “regime change”. Therefore, just like the players in Sobel’s [18] model, the opportunistic type values information about B’s weapons. With full ambiguity, the probability that B invests is decreasing in the cost of investing and increasing in the value of advanced weapons in case of a conflict. If the cost of investing is low and/or the value of advanced weapons is large, then ambiguity does not significantly reduce the risk of weapons proliferation, so the opportunistic type prefers arms inspections in this case. But if the cost is high and/or the value is small, then ambiguity makes A better off regardless of type.

The fact that different types of player A disagree about whether ambiguity is desirable suggests a role for cheap talk. Suppose player A can send either a “tough” or a “conciliatory” message. The conciliatory message allows B to preserve ambiguity about his weapons capabilities, which reduces B’s incentive to invest and thereby reduces the risk of arms proliferation. The tough message, which can be interpreted as insisting that B signs the NPT, removes the ambiguity and thereby increases the risk of arms proliferation. In the communication equilibrium, A uses a “non-convex” strategy. If A is a peaceful “dove” or an aggressive “hawk”, then he has a armed intrinsic preference for a particular action (“don’t attack” for doves, “attack” for hawks). These extreme types do not need arms inspections in order to decide what to do. Since they just want to reduce the risk of arms proliferation, they send the conciliatory message. It is the *intermediate* types who send the tough message. These intermediate types are opportunists who want to

know if B is armed. (To be precise, some opportunists who are “almost” hawks or “almost” doves send the conciliatory message.) We show in the Appendix that this is the only equilibrium where cheap-talk is effective in influencing B’s investment decision. This uniqueness (modulo relabelling the messages) comes from the fact that incentive compatibility requires that a message which makes inspections more likely must increase the risk of arms proliferation.

In the communication equilibrium, player A can trigger inspections by sending the tough message. By revealed preference, all of A’s types prefer the communication equilibrium to the equilibrium with full disclosure. But B may prefer full disclosure if ambiguity is not an effective deterrent. In both the communication equilibrium and the full ambiguity equilibrium, some opportunistic types will attack even though B is armed. The frequency of such “mistakes” determines whether or not ambiguity is good for the small power. There are parameter values where strategic ambiguity about the small power’s arsenal is good for the big power but bad for the small power. The key point is that ambiguity is a substitute for actually acquiring advanced weapons, so more ambiguity means less arms proliferation.

The literature on cheap talk games where the sender has private information was pioneered by Crawford and Sobel [6]. In our model, both the sender (A) and the receiver (B) take actions. Although A’s types are ordered in a natural way according to their propensity to attack, only *intermediate* (opportunistic) types have a *demand for information*. A similar non-convexity appeared in Baliga and Sjöström [1], but it is absent from sender-receiver games of the kind studied by Crawford and Sobel [6] (where the sender takes no action). There is a related literature on financial intermediation and auditing, where costly inspections are used to verify incomes (Townsend [21], Diamond [7], Bond [3], Border and Sobel [4] and Mookherjee and Png [12]). In this literature, the value of information is traded off against the resource cost of inspections. In our model, inspections do not consume significant real resources. We instead focus on a commitment problem: player A cannot commit not to attack if player B reveals that he is unarmed. We show that the optimal policy may be to forego inspections in order to allow B the security to take an action which is good for A (i.e., not to invest).

The paper proceeds as follows. In Section 2, we describe the model. In Section 3, we analyze equilibria without communication, where B either *always* or *never* allows inspections. In Section 4 we consider communication equilibria where A’s message determines whether or not inspections occur.

Section 5 concludes.

## 2 The Arms Proliferation Game

### 2.1 Strategies and Payoffs

There are two players, A and B. Initially, player B has no advanced weapons, but he can try to improve his capabilities by making an investment. His investment decision is binary: he either invests, or he doesn't invest. If B doesn't invest, then he will not acquire advanced weapons. If B invests, he acquires weapons with probability  $\sigma \in (0, 1)$ . The cost of investing is  $k > 0$ . If B acquires advanced weapons then B is *armed*, otherwise B is *unarmed*. Notice that if B invests then he will be armed with probability  $\sigma$ , but if he doesn't invest then he is unarmed for sure. Player A cannot directly observe if B invests or is armed. However, whether B is armed or not is hard information which can be verified by (perfectly reliable) inspectors. If there is an inspection, then B incurs a small cost  $\varepsilon$  which is drawn from a distribution with support  $[0, \bar{\varepsilon}]$  and density  $h$ . The inspection publicly reveals whether B is armed or unarmed.

In the final stage of the game, A decides whether or not to attack B. If A attacks, then A gets a benefit  $a$  and B suffers a cost  $\alpha$ . We refer to  $a$  as player A's *type*. It is A's private information. Player B thinks  $a$  has a continuous distribution with support  $[a_0, a_1]$ , where  $a_0 < 0 < a_1$ . The density is denoted  $f$  and the c.d.f. is denoted  $F$ . If B is armed with advanced weapons, then he can use them if A attacks. This yields an expected benefit  $\gamma \in (0, \alpha)$  for B but an expected cost  $c > 0$  for A. It is useful to define the *normalized cost of investing* to be

$$\kappa \equiv \frac{k}{\sigma\gamma}.$$

Player B has two possible types, “crazy” and “normal”. Player B's type is denoted  $t \in \{z, n\}$  where  $z$  denotes crazy and  $n$  normal. Player A thinks the probability that B is crazy is  $\tau$ . If B is armed but A does not attack, then A's and B's payoffs depend on B's type  $t$ : A suffers a cost  $d_t$  and B derives a benefit  $\delta_t$ . We assume  $0 < d_n < d_z$  and  $\delta_n = 0 < \delta_z < \gamma$ . The crazy type may share the advanced weapons with terrorists, or use them for some other purpose that could hurt A. Therefore, it is more costly for A if a crazy type obtains advanced weapons than if a normal type obtains them:  $d_z > d_n$ . We

assume  $d_n > 0$  because weapons proliferation could be costly to A even if B is not crazy (for example, terrorists may get hold of the technology even if B is normal). The assumption  $\delta_z > \delta_n$  means the advanced weapons are intrinsically more valuable to a crazy type than to a normal type. (We set  $\delta_n = 0$  for simplicity, but we would obtain the same results with  $\delta_n > 0$ .) To simplify, we assume that if A attacks B then he eliminates the threat posed by B. (More generally, the threat could be reduced but not completely eliminated.) While weapons inspectors can verify if B is armed, they cannot verify if his type is crazy or normal. Thus, B's type is his soft information.

The payoffs are summarized in the following matrix.

	B is armed	B is unarmed
A attacks	$a - c, -\alpha + \gamma$	$a, -\alpha$
No attack	$-d_t, \delta_t$	$0, 0$

This payoff matrix does not include B's cost of investment and the cost of inspection. For example, if B invests but does not acquire advanced weapons, there is an inspection, and A attacks, then B's final payoff is  $-\alpha - k - \varepsilon$ .

The solution concept is perfect Bayesian equilibrium. Along the equilibrium path, each player's beliefs are computed from the equilibrium strategies using Bayesian updating. Given these beliefs, each player's behavior must be sequentially rational.

## 2.2 Time Line

The time line is as follows.

Time 0: Player A privately learns  $a \in [a_0, a_1]$  and player B privately learns  $t \in \{z, n\}$ .

Time 1: Cheap-talk stage.

Time 2: Player B decides whether or not to invest.

Time 3: If B invested, then he privately learns whether or not he has acquired weapons.

Time 4: The cost  $\varepsilon$  is realized and then player B decides whether or not to allow inspections. If inspections take place, then the inspectors publicly reveal whether or not B is armed.

Time 5: Player A decides whether or not to attack.

## 2.3 Parameter Restrictions

Player A is said to be a *dove* if  $a < 0$ , an *opportunist* if  $0 < a < c - (\tau d_z + (1 - \tau) d_n)$ , and a *hawk* if  $a > c - (\tau d_z + (1 - \tau) d_n)$ . The probability that A is a dove is  $D \equiv F(0) > 0$ . The probability that he is a hawk is  $H \equiv 1 - F(c - (\tau d_z + (1 - \tau) d_n)) > 0$ .

**Assumption 1:**  $\tau d_z + (1 - \tau) d_n < c < d_z$ .

To interpret this assumption, notice that the first inequality in Assumption 1 implies that

$$a > a - c - (-\tau d_z + (1 - \tau) d_n). \quad (1)$$

The left hand side of (1) is A's gain from attacking an unarmed B. The right hand side of (1) is A's *net* benefit of attacking an armed B if B is crazy with probability  $\tau$  (i.e., the payoff from attacking *minus* the expected payoff from not attacking). Thus, the inequality implies that *given* A's prior  $\tau$ , A is *less* inclined to attack when B is armed than when B is unarmed. Without this inequality, there would be no opportunistic types, advanced weapons could never deter attacks, and the problem would not be interesting. Now the net benefit from attacking an armed B who is thought to be crazy for sure is  $a - c - (-d_z)$ . The second inequality of Assumption 1 says that  $a < a - c - (-d_z)$ . Thus, if A is convinced that B is crazy then A is *more* inclined to attack when B is armed than when B is unarmed. If this inequality were violated, then regardless of A's beliefs, disclosing advanced weapons would always make A less likely to attack. In this case, the unique equilibrium would involve full disclosure, just as in Sobel [18]. To support an equilibrium with ambiguity, A's disutility from allowing a crazy type to get advanced weapons must be sufficiently big. (Nuclear weapons in the hands of North Korea or Iran may not have the range to reach the continental United States, but they do have the potential to destabilize world security in the hands of terrorists.)

Our second assumption guarantees that the cost of inspections is small enough so that it does not significantly influence the set of equilibria.

**Assumption 2:**

$$\bar{\epsilon} < \min\{(F(c - (\tau d_z + (1 - \tau) d_n)) - F(c - d_z))(\alpha - \gamma - \delta_z), (F(0) - F(c - d_z))\alpha, F(0) - F(c$$

Our final assumption ensures that the cost of investing is small enough, so that A cannot achieve his “bliss point”.

**Assumption 3:**

$$\frac{k}{\sigma} < (1 - F(c - d_z))(-\alpha + \gamma) + (1 - F(0))\alpha$$

If Assumption 3 were violated, then the cost of investing would be so high that there would be an equilibrium where a normal type of B never invests and always allows arms inspections. This case would not be very interesting. However, many results, such as the structure and existence of communication equilibrium, in fact do not rely on Assumption 3.

We end this section with a useful preliminary result.

**Proposition 1** *In any perfect Bayesian equilibrium, the crazy type of player B invests with probability one.*

**Proof.** Consider a perfect Bayesian equilibrium. Let  $x(m, t)$  denote the probability that type  $t \in \{z, n\}$  invests after player A has sent message  $m$ . If  $x(m, n) > 0$  for some  $m \in M$  then  $x(m, z) = 1$ . This follows from  $\delta_z > \delta_n$ , which makes the crazy type strictly more willing to invest. Conversely,  $x(m, z) < 1$  implies  $x(m, n) = 0$ .

In order to obtain a contradiction, suppose  $x(m, z) < 1$  for some  $m$ , and let  $M^* \subseteq M$  be the set of messages that minimize  $x(m, z)$ . If  $m^* \in M$  then  $x(m^*, z) < 1$  and  $x(m^*, n) = 0$ .

*Claim 1:* If  $m^* \in M^*$  then  $0 < x(m^*, z) < 1$  for .

Proof of claim: By hypothesis,  $x(m^*, z) < 1$ . Suppose  $x(m^*, z) = 0$ . In this case, B will be unarmed for sure following message  $m^*$ . Clearly, all of A’s types will send a message in  $M^*$ , since a zero probability that B invests is the best possible outcome for A. Since B will be known to be unarmed, player A attacks if and only if  $a \geq 0$ , which happens with probability  $1 - F(0)$ . Suppose B changes to a strategy where he invests, and refuses inspections if the investment succeeds. Since A never attacks when  $a - c < -d_z$ , the probability of an attack can be at most  $1 - F(c - d_z)$ , and B’s expected improvement will be at least

$$\sigma \{(1 - F(c - d_z))(-\alpha + \gamma)) + F(c - d_z)\delta_t - (1 - F(0))(-\alpha)\} - k$$

This is strictly positive by Assumption 3. This contradiction proves the claim.

*Claim 2:* Player A must send a message in  $M^*$  if  $a > 0$  or  $a < c - d_z$ .

Proof of claim: If  $a - c > -d_n$  then A always attacks, so at the message stage he simply wants to minimize the probability that B is armed. He does this by sending a message in  $M^*$ . Similarly, if  $a - c < -d_z$  then A never attacks, and again he wants to minimize the probability that B is armed. Finally, consider the case  $0 < a \leq c - d_n$ . If A sends  $m^* \in M^*$  and then (regardless of what happens at the inspections stage) attacks for sure, then his expected payoff is  $a - \sigma\tau x(m^*, z)c$ . Suppose instead he sends  $m' \notin M^*$ . Following this message, B's crazy type will be armed with probability  $\sigma x(m', z)$  and his normal type will be armed with probability  $\sigma x(m', n)$ . If B is unarmed, A prefers to attack (since  $a > 0$ ). If B is armed, then A prefers to attack if and only if B is crazy (since  $-d_z < a - c < -d_n$ ). Therefore, A's maximum possible payoff from sending message  $m'$  is

$$\begin{aligned} & (1 - \tau\sigma x(m', z) - (1 - \tau)\sigma x(m', n))a + \tau\sigma x(m', z)(a - c) + (1 - \tau)\sigma x(m', n)(-d_n) \\ = & (1 - \sigma(1 - \tau)x(m', n))a - \sigma\tau x(m', z)c - \sigma(1 - \tau)x(m', n)d_n \\ < & a - \sigma\tau x(m^*, z)c \end{aligned}$$

since  $a > 0$  and  $x(m^*, z) < x(m', z)$ . This proves the claim.

Notice that claim 2 implies  $M^* \neq \emptyset$ .

*Claim 3:* If any  $m^* \in M^*$  was sent and B is armed, then B will not allow inspections.

Proof of claim: To obtain a contradiction, suppose that following  $m^* \in M^*$ , there is a positive probability that inspections are allowed and reveal that B is armed. Since only crazy types invest ( $x(m^*, z) > 0 = x(m^*, n)$ ), B will be known to be crazy once he reveals that he is armed. Thus, all types of A with  $a \geq c - d_z$  will attack. Since type  $a < c - d_z$  will never attack in any situation, if B is armed he is strictly better off not revealing it, since inspections are costly and will not reduce the probability of attack. This contradiction proves the claim.

*Claim 4:* Following any  $m^* \in M^*$ , there must be positive probability that B is unarmed and refuses inspections.

Proof of claim: Suppose that if B observes  $m^* \in M^*$  and is unarmed, then he allows inspections with probability one. It follows that if B refuses inspections, he will be known to be armed. Therefore, following message  $m^*$  there is never ambiguity about whether or not B is armed. Since  $m^* \in M^*$

both minimizes the probability that B arms and also fully reveals if B is armed or not, all types of A will either send  $m^*$ , or else some equivalent message in  $M^*$  (i.e., a message which also fully reveals if B is armed). Since  $x(m^*, n) = 0$  for  $m^* \in M^*$ , the normal type will never invest in this equilibrium, will always allow inspections, and will be attacked whenever  $a > 0$ . So his expected payoff is

$$(1 - F(0))(-\alpha) - E\varepsilon \quad (2)$$

where  $E\varepsilon$  is the expected value of  $\varepsilon$ . Now, in no circumstance would A ever attack if  $a < c - d_z$ . Therefore, if B invests and then allows inspections if and only if the investment fails, his expected payoff is at least

$$\sigma \{(1 - F(c - d_z))(-\alpha + \gamma) + F(c - d_z)\delta_t\} + (1 - \sigma) \{(1 - F(0))(-\alpha) - E\varepsilon\} - k \quad (3)$$

Assumption 3 implies that (3) is strictly greater than (2) for both types of B. Thus, the normal type can improve on his equilibrium payoff by investing. This contradiction proves the claim.

Now we can complete the proof of the Proposition. Let  $\zeta$  denote the probability that A attacks, conditional on a message in  $M^*$  having been sent and inspections refused. For any  $m^* \in M^*$ ,  $x(m^*, z) > 0 = x(m^*, n)$ , so following message  $m^*$  player B will be either unarmed or armed and crazy. In either case, A prefers to attack if  $a \geq 0$ . Thus, by claim 2, if  $a \geq 0$  then A sends a message in  $M^*$  and then attacks for sure. Type  $a < c - d_z$  also sends a message  $M^*$ , but never attacks. Therefore,  $1 - F(0) \leq \zeta \leq 1 - F(c - d_z)$ . Following any  $m^* \in M^*$ , when B is unarmed he at least weakly prefers to refuse inspections. Thus, conditional on some message in  $M^*$  having been sent and player B not investing, B's expected payoff is  $-\zeta\alpha$ . On the other hand, if after any  $m^* \in M^*$  B invests and then refuses inspections, then his expected payoff conditional on a message in  $M^*$  is

$$\begin{aligned} & \sigma \{\zeta(-\alpha + \gamma) + (1 - \zeta)\delta_t\} + (1 - \sigma)\zeta(-\alpha) - k \\ \geq & \sigma(1 - F(c - d_z))(-\alpha + \gamma) + (1 - \sigma)\zeta(-\alpha) - k \\ > & -\sigma(1 - F(0))\alpha - (1 - \sigma)\zeta\alpha \geq -\zeta\alpha \end{aligned}$$

where the first inequality is due to  $\zeta \leq 1 - F(c - d_z)$  and  $\delta_t \geq 0$ , the second to Assumption 3, and the third to  $1 - F(0) \leq \zeta$ . Therefore, both types of B will strictly prefer to invest in response conditional on receiving a message in  $M^*$ , a contradiction of  $x(m^*, n) = 0$ . ■

## 3 Equilibria Without Communication

### 3.1 Equilibria with Full Disclosure

In an equilibrium with *full disclosure*, there is never any ambiguity about B's weapons on the equilibrium path. For example, an inspection may occur with probability one, or B may allow inspections if and only if he is armed (in which case a refusal to allow inspections reveals that B is unarmed). With full disclosure, all of A's types will send a message which minimizes the probability that B invests. Thus, cheap-talk cannot be effective, i.e., the probability that B invests cannot depend on A's type. Without ambiguity, communication cannot prevent arms proliferation. (Since B has no incentive to reveal that he is crazy, we can without loss of generality assume that B sends no message.)

**Proposition 2** *There is an equilibrium with full disclosure. Full disclosure implies that both types of player B invest with probability one. Cheap-talk cannot reduce the probability that B invests.*

**Proof.** Suppose in equilibrium, B refrains from investing. With full disclosure, an unarmed B is attacked whenever  $a \geq 0$ , which happens with probability  $1 - F(0)$ . Consider a deviation where B invests and refuses inspections if successful, which triggers an attack with *at most* probability  $1 - F(c - d_z)$  (because A will never attack if  $a < c - d_z$ ). The gain from this deviation is *at least*

$$\sigma \{ (1 - F(c - d_z)) (-\alpha + \gamma) + \sigma F(c - d_z) \delta_t - (1 - F(0)) (-\alpha) \} - k$$

This expression is strictly positive, by Assumption 3. Therefore, in any full disclosure equilibrium, player B invests with probability one.

It remains to show that an equilibrium with full disclosure exists. Let the equilibrium strategy specify that player B invests with probability one, and he allows inspections if and only if he is armed. If inspections reveal that B is armed, then A attacks if  $a - c > \tau d_z + (1 - \tau) d_n$  (A thinks B is crazy with probability  $\tau$  since both types are armed with probability  $\sigma$ ). If B refuses inspections, then A infers that B is unarmed, so A attacks if  $a \geq 0$ . If B should allow inspections even though he is unarmed, he is still attacked if  $a \geq 0$ , so he has no reason to allow inspections in this case. Suppose B

deviates by refusing inspections when he is armed. This raises the probability of attack from  $1 - F(c - (\tau d_z + (1 - \tau) d_n))$  to  $1 - F(0)$ , which has a cost

$$(F(c - (\tau d_z + (1 - \tau) d_n)) - F(0)) (\alpha - \gamma - \delta_t)$$

The gain from the deviation is only  $\varepsilon$ . Assumption 2 guarantees that the cost exceeds the benefit, so B prefers to reveal if he is armed. Finally, given full disclosure, Assumption 3 guarantees that B prefers to invest. ■

### 3.2 Equilibria with Full Ambiguity

Proposition 2 implies that B prefers to invest unless there is some ambiguity about his capabilities. In an equilibrium with *full ambiguity*, inspections never occur on the equilibrium path. Clearly, with full ambiguity communication cannot be effective. For ambiguity to deter A from attacking, the normal type of B must invest with sufficiently high probability. However, this requires that A attacks with sufficiently high probability, or else the normal type has no incentive to invest. The required equilibrium probabilities will depend on the normalized cost of investing,  $\kappa$ .

**Proposition 3** *There is an equilibrium with full ambiguity. Full ambiguity implies that B's normal type invests with probability  $\tilde{x}$ , where  $0 < \tilde{x} < 1$  if*

$$\kappa > 1 - F(\sigma(c - \tau d_z - (1 - \tau) d_n)), \quad (4)$$

*and  $\tilde{x} = 1$  otherwise. Cheap-talk cannot reduce the probability that B invests.*

**Proof.** Proposition 1 implies that player B's crazy type always invests. If B never allows inspections, then all of A's types simply want to minimize the probability that B's normal type invests. Therefore, the probability of investment must be independent of A's type. Let  $\tilde{x}$  denote the probability that the normal type of B invests. Thus, B is armed with probability  $\sigma(\tau + (1 - \tau)\tilde{x})$ .

The equilibrium must satisfy a cut-off property: there is  $\tilde{a}$  such that if there is no inspection then A attacks if and only if  $a > \tilde{a}$ . In equilibrium, doves will not attack but hawks will. Type  $\tilde{a} \in (a_0, a_1)$  must be indifferent between attacking and not attacking. Type  $\tilde{a}$  expects  $\tilde{a} - \sigma(\tau + (1 - \tau)\tilde{x})c$

by attacking, and  $-\sigma(\tau d_z + (1 - \tau)\tilde{x}d_n)$  by not attacking. The cut-off property implies that type  $\tilde{a}$  must be indifferent between his two actions. This indifference condition holds if

$$\tilde{a} = \sigma\tau(c - d_z) + \sigma(1 - \tau)\tilde{x}(c - d_n) \quad (5)$$

If B deviates by allowing inspections, and he is found to be unarmed, then A attacks if and only if  $a > 0$ . But if B is found to be armed, then we may suppose A attacks if and only if  $a > c - d_z$ . This is supported by the off-the-equilibrium path belief that B is crazy (which is the belief most likely to support the equilibrium, since it punishes B's deviation most strictly). We will show that  $\tilde{a} > 0$ , so inspections always increase the probability of attack. This clearly implies that B prefers to refuse inspections.

If  $0 < \tilde{x} < 1$  then B's normal type must be indifferent between investing and not investing. Since B is attacked with probability  $1 - F(\tilde{a})$ , he is indifferent between investing and not investing if

$$-(1 - F(\tilde{a}))\alpha = -(1 - F(\tilde{a}))(\alpha - \sigma\gamma) - k \quad (6)$$

which is the same as

$$\kappa - (1 - F(\tilde{a})) = 0 \quad (7)$$

If  $\kappa > 1 - F(\tilde{a})$  then the normal type's unique best response is not to invest, so  $\tilde{x} = 0$ . Similarly, if  $\kappa < 1 - F(\tilde{a})$  then the unique best response implies  $\tilde{x} = 1$ . Define

$$\Gamma(x) \equiv \kappa - (1 - F(\sigma\tau(c - d_z) + \sigma(1 - \tau)x(c - d_n)))$$

An equilibrium where  $0 < \tilde{x} < 1$  requires that both (7) and (5) hold, which implies  $\Gamma(\tilde{x}) = 0$ . Now,

$$\frac{k}{\sigma} < (F(0) - F(c - d_z))(-\alpha) + (1 - F(c - d_z))\gamma < (1 - F(0))\gamma < (1 - F(\sigma\tau(c - d_z)))\gamma \quad (8)$$

which implies  $\kappa < 1 - F(\sigma\tau(c - d_z))$ . (The first inequality in (8) follows from Assumption 3, the second follows from  $c - d_z < 0$  and  $\alpha > \gamma$ , the third follows from  $\sigma\tau(c - d_z) < 0$ .) Therefore,  $\Gamma(0) < 0$ . Since  $\Gamma'(x) > 0$ , there is  $\tilde{x} \in (0, 1)$  such that  $\Gamma(\tilde{x}) = 0$  if and only if  $\Gamma(1) > 0$ , which is equivalent to (4). Thus, there are two possible cases.

*Case (i):* (4) holds. In this case, then there is  $\tilde{x} \in (0, 1)$  such that  $\Gamma(\tilde{x}) = 0$ , and this is the only candidate for a full ambiguity equilibrium.

(Since  $\Gamma(0) < 0 < \Gamma(1)$ , it is not possible that the normal type invests with probability 0 or 1). It is indeed an equilibrium, because  $\tilde{a} > 0$ . This follows from

$$\kappa - (1 - F(0)) < 0 = \kappa - (1 - F(\tilde{a})).$$

where the first inequality is due to (8) and the second to (7).

*Case (ii):* (4) is violated. In this case,  $\Gamma(1) \leq 0$  so we must have  $\tilde{x} = 1$ . It is indeed an equilibrium, because (5) and Assumption 1 yield

$$\tilde{a} = \sigma(c - \tau d_z - (1 - \tau)d_n) > 0.$$

■

**Remark.** In the proof of Proposition 3, equilibrium is supported by the belief that B is crazy if inspections revealed that he is armed. These out-of-equilibrium beliefs qualify as reasonable according to standard arguments, such as the D1 criterion of Banks and Sobel [2]. Indeed, the equilibrium payoff of an armed normal type at the inspection stage is  $(-\alpha + \gamma)(1 - F(\tilde{a}))$ . If he allows inspections, and A attacks with probability  $p$ , then his expected payoff is  $(-\alpha + \gamma)p - \varepsilon$ . Thus, inspections would make the armed normal type weakly better off if and only if

$$(-\alpha + \gamma)p - \varepsilon \geq (-\alpha + \gamma)(1 - F(\tilde{a})). \quad (9)$$

Similarly, inspections would make the armed crazy type strictly better off if and only if

$$(-\alpha + \gamma)p + (1 - p)\delta_z - \varepsilon > (-\alpha + \gamma)(1 - F(\tilde{a})). \quad (10)$$

Since (9) implies (10), the D1 criterion suggests that the crazy type should be assigned probability one if out-of-equilibrium inspections reveal he is armed.

In equilibrium, ambiguity has its price, because some opportunistic types attack even though B is armed. The welfare implications of ambiguity depend on the probability of such “mistakes”. It is useful to define

$$a^* \equiv \frac{\sigma d_n (c - \tau d_z - (1 - \tau) d_n)}{(1 - \sigma)c + \sigma d_n} \quad (11)$$

Assumption 1 implies

$$0 < a^* < \sigma (c - (\tau d_z + (1 - \tau) d_n)), \quad (12)$$

so type  $a^*$  is an opportunist. To look at the implications of ambiguity, we can distinguish two cases.

*Case 1:* Suppose the normalized cost of developing advanced weapons is high:

$$\kappa > 1 - F(a^*), \quad (13)$$

where  $a^*$  is defined by (11). The inequalities (12) and (13) imply that (4) holds, so B refrains from investing with probability  $1 - \tilde{x} > 0$  behind the veil of ambiguity. Clearly, hawks and doves strictly prefer full ambiguity to full disclosure (under full disclosure B invests with probability one). Among the opportunists, it is not hard to see that the one most likely to want inspections is precisely type  $\tilde{a}$ , defined by . The smaller is  $\tilde{x}$ , the more likely it is that type  $\tilde{a}$  prefers full ambiguity. Type  $\tilde{a}$ 's expected utility under full ambiguity is  $\tilde{a} - \sigma(\tau + (1 - \tau)\tilde{x})c$ . Compare this to the outcome where B always invests and allows inspections. After the inspection, type  $\tilde{a}$  attacks if and only if B is unarmed, which happens with probability  $1 - \sigma$  (from (5), type  $\tilde{a}$  is an opportunist). Thus, type  $\tilde{a}$ 's expected payoff would be

$$(1 - \sigma)\tilde{a} - \sigma(\tau d_z + (1 - \tau)d_n).$$

Type  $\tilde{a}$  prefers full ambiguity if and only if

$$(1 - \sigma)\tilde{a} - \sigma(\tau d_z + (1 - \tau)d_n) < \tilde{a} - \sigma(\tau + (1 - \tau)\tilde{x})c \quad (14)$$

Using the definition of  $\tilde{a}$ , (14) is equivalent to  $\tilde{x} < x^*$ , where

$$x^* \equiv \frac{(1 - \sigma)\tau(d_z - c) + (1 - \tau)d_n}{(1 - \sigma)(1 - \tau)c + \sigma(1 - \tau)d_n},$$

The first inequality in Assumption 1 implies  $x^* < 1$ . Clearly,  $\tilde{x} < x^*$  if  $\tilde{x} = 0$ . Suppose instead that  $\tilde{x} > 0$ . Since  $\Gamma(0) < 0 = \Gamma(\tilde{x}) < \Gamma(1)$  and  $\Gamma'(x) > 0$ , we have  $\tilde{x} < x^*$  if and only if  $\Gamma(x^*) > 0$ , which is equivalent to (13). Thus, in case 1 ambiguity reduces the risk of arms proliferation sufficiently to make all of A's types better off.

*Case 2:* Suppose the normalized cost of developing advanced weapons is low:

$$\kappa < 1 - F(a^*). \quad (15)$$

If (4) holds, then B invests with probability  $\tilde{x} < 1$ . Therefore, hawks and doves strictly prefer full ambiguity to full disclosure. However, by a

similar reasoning as in case 1, we find that (15) implies that type  $\tilde{a}$  strictly prefers full disclosure to full ambiguity (inequality (14) is reversed). If (4) is violated, then B invests with probability one under full ambiguity, so some opportunistic types strictly prefer full disclosure, because it allows them to make better decisions. In case 2, ambiguity does not significantly reduce the risk of arms proliferation. Therefore, there are always opportunistic types of A who prefer disclosure.

So far, we have considered only A's welfare. Now consider the situation from the point of view of B. With full ambiguity, player A attacks when  $a \geq \tilde{a}$ . With full disclosure, player A attacks if  $a \geq c - \tau d_z - (1 - \tau) d_n$  when B is armed, and if  $a \geq 0$  when B is unarmed. Therefore, when moving from full ambiguity to full disclosure, player B's expected utility (not including the cost of inspection) changes by an amount

$$\sigma(\alpha - \gamma - \delta_t) [F(c - \tau d_z - (1 - \tau) d_n) - F(\tilde{a})] - (1 - \sigma) \alpha [F(\tilde{a}) - F(0)]. \quad (16)$$

The first term is positive. This term is due to the fact that with full ambiguity, a measure  $F(c - \tau d_z - (1 - \tau) d_n) - F(\tilde{a})$  of "tough" opportunists attack B when he is armed. (Under full disclosure, B's weapons would be revealed and the tough opportunists would be deterred.) The second term is negative. It is due to the fact that with full ambiguity, a measure  $F(\tilde{a}) - F(0)$  of "weak" opportunists do not attack B when he is unarmed. (Under full disclosure, the weak opportunists would attack the unarmed B.) Thus, disclosure deters "tough" opportunists when B is armed, but ambiguity deters "weak" opportunists when B is unarmed. Without making further assumptions on the distribution of A's types we cannot sign the expression in (16).

We summarize these findings in the following proposition:

**Proposition 4** *All of A's types prefer full ambiguity to full disclosure if and only if (13) holds. Player B prefers full ambiguity to full disclosure if and only if the expression in (16) is negative.*

Clearly, hawks and doves always at least weakly prefer full ambiguity to full disclosure, since they have nothing to gain from inspections (their actions will not depend on the arms inspector's report). However, there is a case, namely if (4) holds in Case 2, where hawks and doves *strictly* prefer ambiguity while some opportunistic types *strictly* prefer disclosure. In this case there is a conflict of interest among A's types about whether inspections

are desirable or not. In order to alleviate this conflict of interest, B's decision to allow inspections should be made to depend on A's type.

## 4 Communication Equilibrium

Inspections generate information about B's capabilities. This information may be more or less useful, depending on A's type. The extreme types (hawks and doves) dislike arms proliferation, but they do not benefit from inspections *per se*, because they will not act on the information that is generated. The intermediate types are opportunistic and their optimal action depends on whether B is armed or not. Hence, they are willing to trade-off a higher probability of arms proliferation for a higher probability of inspections. This allows us to construct an equilibrium with two informative messages. The intermediate types send a "tough" message that leads to inspections, but also induces B to invest. The extreme types send a "conciliatory" message which does not lead to inspections, but also reduces the risk that B will invest.

The communication equilibrium exists if and only if two conditions are satisfied. First, we must be in case 2 of Section 3.2. Otherwise, all of A's types would prefer ambiguity, and no-one would send the "tough" message. Thus, the first condition is that (15) must hold. Second, the prior probability that A is a hawk must be small. Otherwise, B would invest for sure after the conciliatory message, which is sent by both hawks and doves. In any case, it is intuitively clear that if A is likely to be a hawk then ambiguity will not prevent B from investing. Specifically, the second condition turns out to be

$$\frac{H}{H + D} < \kappa. \quad (17)$$

**Proposition 5** *Suppose*

$$\frac{H}{H + D} < \kappa < 1 - F(a^*) \quad (18)$$

*There is a communication equilibrium where, for some  $a'$  and  $a''$ , player A sends a "tough" message if  $a \in (a', a'')$  and a "conciliatory" message otherwise. Player B is more likely to invest and more likely to allow inspections following the tough message.*

**Proof.** Consider the following strategies. There is  $a'$  and  $a''$ , where  $0 < a' < a'' < c - \tau d_z - (1 - \tau) d_n$ , such that A sends the tough message if  $a \in (a', a'')$  and the conciliatory message otherwise. Player B allows inspections if and only if he hears the tough message and is armed. Player B's crazy type invests with probability one, regardless of message.

If A sends the tough message, then B's normal type invests with probability 1, and allow inspections if and only if the investment is successful. Since type  $a \in (a', a'')$  who sends the tough message is opportunistic, he will attack if inspections are refused or if inspections reveal that B is unarmed. If they reveal that B is armed, then he will not attack.

If A sends the conciliatory message, then the normal type of B invests with probability  $x \in (0, 1)$  and refuses inspections. If there is no inspection following the conciliatory message, then A attacks if  $a \geq a''$  but not if  $a \leq a'$ . If there is a "surprise inspection" which reveals that B is armed, then A attacks if and only if  $a > c - d_z$ . This is justified by the out-of-equilibrium belief that B is a crazy type (which is consistent with the D1 criterion). If the surprise inspection reveals that B is unarmed, then A attacks if and only  $a > 0$ .

We now make sure that  $a'$  and  $a''$  are indifferent between the two messages. Suppose type  $a''$  sends the tough message. With probability  $\sigma$ , player B is armed and type  $a''$  gets  $-\tau d_z - (1 - \tau) d_n > a'' - c$ . With probability  $1 - \sigma$ , player B is unarmed and type  $a''$  gets  $a'' > 0$ . Thus, the expected payoff is  $(1 - \sigma)a'' - \sigma\tau d_z - \sigma(1 - \tau) d_n$ . If type  $a''$  sends the conciliatory message then B will be armed with probability  $\sigma(\tau + (1 - \tau)x)$ . Type  $a''$  will attack, and get expected payoff  $a'' - \sigma(\tau + (1 - \tau)x)c$  (we verify later that attacking is optimal). For type  $a''$  to be indifferent between the two messages, we must have

$$a'' = (\tau + (1 - \tau)x)c - \tau d_z - (1 - \tau) d_n < c - \tau d_z - (1 - \tau) d_n \quad (19)$$

If type  $a'$  sends the tough message, his expected payoff is  $(1 - \sigma)a' - \sigma\tau d_z - \sigma(1 - \tau) d_n$ . If type  $a'$  sends the conciliatory message, he will not attack (we verify later that this is optimal), and he gets expected payoff  $-\sigma\tau d_z - \sigma(1 - \tau)x d_n$ . For type  $a'$  to be indifferent, we must have

$$a' = \frac{(1 - x)\sigma(1 - \tau) d_n}{1 - \sigma} > 0 \quad (20)$$

Define

$$x^* \equiv \frac{(1 - \sigma)\tau(d_z - c) + (1 - \tau) d_n}{(1 - \sigma)(1 - \tau)c + \sigma(1 - \tau) d_n} < 1 \quad (21)$$

where the inequality follows from the first inequality in Assumption 1. If  $x = x^*$  is substituted into (19) and (20), we get  $a' = a'' = a^*$ , as defined in (11). Now (19) and (20) imply that  $a''$  is increasing in  $x$  and  $a'$  is decreasing in  $x$ . Thus,  $a'' > a'$  as long as  $x > x^*$ .

We now verify B's incentive to play according to his strategy. First, consider the decision to allow inspections. If he hears the tough message but is unarmed, then B realizes he will be attacked whether or not he allows inspections. He strictly prefers to refuse inspections to save the cost  $\varepsilon$ . If B is armed then his expected payoff from allowing inspections following the tough message is  $\delta_t - \varepsilon$ , while his expected payoff from refusing is  $-(\alpha - \gamma)$ . He prefers to allow inspections as  $\delta_t - \varepsilon > -(\alpha - \gamma)$  by Assumption 2. Similarly, B strictly prefers to refuse inspections after the conciliatory message since inspections only increase the probability of an attack.

Next, consider the normal type's decision to invest. If B hears the tough message, then his expected payoff from investing is  $\sigma(-\varepsilon) + (1 - \sigma)(-\alpha)$ . His expected payoff from not investing is  $-\alpha$ . Since  $\alpha > \varepsilon$ , he prefers to invest.

Now consider the normal type's investment decision following the conciliatory message. If B hears the conciliatory message, then he thinks A will attack if  $a \geq a''$  but not if  $a \leq a'$ . Accordingly, if B invests his expected payoff is

$$-\frac{1 - F(a'')}{F(a') + 1 - F(a'')} (\alpha - \sigma\gamma) - k$$

If he does not invest, his expected payoff is

$$-\frac{1 - F(a'')}{F(a') + 1 - F(a'')} \alpha$$

Player B's normal type must be indifferent between investing and not investing (since  $0 < x < 1$ ), which is true if

$$(1 - F(a'') + F(a')) \kappa - (1 - F(a'')) = 0 \tag{22}$$

We can use (19) and (20) to substitute for  $a'$  and  $a''$  in (22). Define

$$\Psi(x) \equiv \left( 1 - F((\tau + (1 - \tau)x)c - \tau d_z - (1 - \tau)d_n) + F\left(\frac{(1 - x)\sigma(1 - \tau)d_n}{1 - \sigma}\right) \right) \kappa - (1 - F((\tau + (23)$$

Notice that

$$\Psi(x^*) = \kappa - (1 - F(a^*))$$

and

$$\Psi(1) = (1 - F(c - \tau d_z - (1 - \tau) d_n) + F(0))\kappa - (1 - F(c - \tau d_z - (1 - \tau) d_n))$$

The indifference condition (22) is verified, together with (19) and (20), if  $x$  is chosen such that  $\Psi(x) = 0$ . Now, (18) is equivalent to  $\Psi(x^*) < 0 < \Psi(1)$ . By continuity, there is  $x \in (x^*, 1)$  such that  $\Psi(x) = 0$ . Since  $x > x^*$ ,  $a'' > a'$ .

Notice that A's extreme types ( $a < a'$  and  $a > a''$ ) are less interested in inspections than the intermediate types. Since types  $a'$  and  $a''$  are indifferent between the two messages, it is indeed optimal for the intermediate types to send the tough message, and for the extreme types to send the conciliatory message. Also, since B's normal type always weakly prefers to invest, it is optimal for the crazy type to always invest.

It remains to verify two assertions made above. First, it should not be optimal for type  $a'$  to send a conciliatory message and then attack. If type  $a'$  chooses such a strategy, then he gets

$$a' - \sigma(\tau + (1 - \tau)x)c = a' - \sigma(a'' + \sigma\tau d_z + \sigma(1 - \tau)d_n) < (1 - \sigma)a' - \sigma(\tau d_z + (1 - \tau)d_n)$$

where the equality uses (19), and the inequality is due to  $a'' > a'$ . The right hand side expression is what type  $a'$  gets in equilibrium.

Second, it should not be optimal for type  $a''$  to send a conciliatory message and then not attack. If type  $a''$  chooses such a strategy, then he gets

$$-\sigma\tau d_z - \sigma(1 - \tau)x d_n = -\sigma\tau d_z - \sigma(1 - \tau)d_n + (1 - \sigma)a' < -\sigma\tau d_z - \sigma(1 - \tau)d_n + (1 - \sigma)a''$$

where the equality uses (20), and the inequality is due to  $a'' > a'$ . The right hand side expression is what type  $a''$  gets in equilibrium. ■

The communication equilibrium has the same “non-convex” structure as the cheap-talk equilibrium in Baliga and Sjöström [1]. Different types trade off “coordination” and “cooperation” at different rates. The intermediate types put a high value on coordination: they need information in order to make an optimal decision. The extreme types mainly want the opponent to cooperate (by not investing). Therefore, it is possible to separate the extreme types from the intermediate types. If the prior probability that A is a hawk is low enough, specifically if  $H/(H + D) < \kappa$ , then the conciliatory message reduces B's fear and lowers the risk of arms proliferation. However, in Baliga and Sjöström [1], all decisions were made *simultaneously* so the issue of strategic ambiguity did not arise. Since the current model allows the

possibility of arms inspections taking place before A decides to attack, we can study the costs and benefits of strategic ambiguity. As we have shown, strategic ambiguity is required for effective communication.

By revealed preference, all of A's types weakly prefer the communication equilibrium to full disclosure (and there is strict preference for some). Consider B's payoff. With full disclosure, if B's investment succeeds then he is attacked with probability  $1 - F(c - \tau d_z - (1 - \tau) d_n)$ , otherwise he is attacked with probability  $1 - F(0)$ . In the communication equilibrium, if A sends the tough message then B is attacked if and only if the investment fails. If A sends the conciliatory message, then B is attacked with probability

$$\frac{1 - F(a'')}{F(a') + 1 - F(a'')}.$$

This implies that if we move from the communications equilibrium to full disclosure, B's expected payoff changes by

$$\sigma(\alpha - \gamma + \delta_t) [F(c - \tau d_z - (1 - \tau) d_n) - F(a'')] - (1 - \sigma) [F(a') - F(0)] \alpha - [1 - (F(a'') - F(a'))] E \quad (24)$$

The interpretation is similar to (16). The first term is positive and is due to the fact that there is a measure  $F(c - \tau d_z - (1 - \tau) d_n) - F(a'')$  of "tough" opportunists who send the conciliatory message but then attack B, even though B may be armed. (Under full disclosure, the tough opportunists are deterred by B's weapons.) The second term is negative and is due to the fact that there is a measure  $F(a') - F(0)$  of "weak" opportunists who send a conciliatory message and then do not attack, even though B may be unarmed. (Under full disclosure, the weak opportunists attack whenever B is unarmed.) Again, disclosure deters "tough" opportunists when B is armed, but ambiguity deters "weak" opportunists when B is unarmed. Without making further assumptions on the distribution  $F$ , we cannot sign the expression in (24).

**Proposition 6** *All of A's types prefer the communication equilibrium to full disclosure. Player B prefers the communication equilibrium to full disclosure if and only if the expression (24) is negative.*

So far, we have considered communication equilibria with just two messages, but this turns out to be without loss of generality.

**Proposition 7** *All equilibria with effective cheap-talk can be replicated by using just two messages..*

We sketch the proof of Proposition 7 here (the rigorous proof is in the Appendix). Let  $M$  be some arbitrary message space, and consider an equilibrium of the game. Let  $M^C$  be the set of “conciliatory” messages that minimize the probability that B invests. Let  $M^T = M \setminus M^C$  be the set of “tough” messages. Any type of A who either *always* or *never* attacks in equilibrium must send a message in  $M^C$ , since all he cares about is reducing the probability that B invests. If cheap talk is effective in influencing B’s decision to invest, then some types must send a message in  $M^T$ , and they must all be opportunists. They do not attack if inspections reveal B is armed, but attack otherwise. Therefore, any message  $m^T \in M^T$  will cause B to invest and reveal his weapons if he is successful. As all messages in  $M^T$  lead to the same outcome, we can assume  $M^T$  is a singleton. Furthermore, B must invest with positive probability in response to any message  $m^C \in M^C$ , and must refuse inspections with positive probability when the investment succeeds (otherwise, all types would prefer to send a message in  $M^C$ ). The proof is finished by showing that, in fact, all messages in  $M^C$  must cause B to invest with the same probability and refuse inspections always. Hence we can assume  $M^C$  is also a singleton.

## 5 Conclusion

In policy debates, it is often argued that U.S. policy should be to eliminate ambiguity (e.g., Schrage [17]). Sobel [18] pointed out that ambiguity makes it more difficult to distinguish the unarmed from the armed, which protects the unarmed from being attacked. This suggests that ambiguity may be good for small powers (whose capabilities are uncertain) but not for big powers (who are known to be armed). However, we argue that once the small power’s incentive to arm itself is taken into account, the opposite may be true. For ambiguity to be part of an equilibrium, the small power (B) must have an incentive to invest with positive probability, which means attacks must be sufficiently likely. Some of these attacks will be “mistakes”: the leader of the big power (A) attacks even though he would have been deterred, had he known the small power’s true capabilities. If such mistakes are very likely, then strategic ambiguity hurts the small power. We stress instead another positive aspect of ambiguity: the small power’s incentive to arm itself is reduced by ambiguity. Therefore, strategic ambiguity tends to benefit the big power, at least if the leader is a type who will not make any “mistakes”.

If A is an opportunistic type, then he needs information about B's true capabilities in order to avoid making mistakes. In a communication equilibrium, opportunistic types send "tough" messages, which can be interpreted as a demand to sign the NPT. Player B responds by revealing his true capabilities. Dovish types instead send "conciliatory" messages. If B hears a conciliatory message, then he maintains a policy of ambiguity, but he is less likely to actually acquire advanced weapons. Unfortunately, hawks have an incentive to masquerade as doves and send a conciliatory message as well. Therefore, the nature of the equilibrium set depends on the relative likelihood of hawks and doves,  $H/(H + D)$ . If  $H/(H + D)$  is too large then a conciliatory message will not reassure B, who suspects a "false dove", and communication will be ineffective.

A second determinant of the equilibrium set is the normalized cost of investing,  $\kappa$ . Recall that  $\kappa$  is higher the bigger is the cost  $k$  of investing; the smaller is the probability  $\sigma$  that B will acquire advanced weapons; and the smaller is the value of advanced weapons,  $\gamma$ . With ambiguity, the probability that B invests is decreasing in  $\kappa$ . If  $\kappa$  is small then B very likely will attempt to get advanced weapons, whether there is ambiguity or not. Thus, the smaller is  $\kappa$ , the more likely it is that the opportunistic type prefers inspections. But if  $\kappa$  is high then ambiguity makes A better off regardless of type.

## References

- [1] Baliga, S. and T. Sjöström (2004), "Arms Races and Negotiations," *Review of Economic Studies* **71**:351-369.
- [2] Banks, J. and J. Sobel (1987), "Equilibrium Selection in Signaling Games," *Econometrica* **55**: 647-662.
- [3] Bond, P. (2004): "Bank and Nonbank Financial Intermediation," *Journal of Finance* **59**: 2489-2530.
- [4] Border, K. and J. Sobel (1987): "Samurai Accountant: A Theory of Auditing and Plunder," *Review of Economic Studies* **54**: 525-540.
- [5] Cohen, A. (1998) *Israel and the Bomb* (New York City: Columbia University Press)

- [6] Crawford, V. and J. Sobel (1982), "Strategic information transmission," *Econometrica* **50**: 1431-1451.
- [7] Diamond, D. (1984), "Financial Intermediation and Delegated Monitoring," *Review of Economic Studies*, **51**: 393-414.
- [8] Grossman, S. (1981), "The Informational Role of Warranties and Private Disclosure about Product Quality", *Journal of Law and Economics* **24**: 461-483
- [9] Hecker, S. (2004), "Visit to the Yonugbyon Nuclear Scientific Research Center in North Korea", written statement to the Senate Committee on Foreign Relations
- [10] Kreps, D. and R. Wilson (1982), "Sequential Equilibrium," *Econometrica* **50**: 863-894
- [11] Milgrom, P. (1981), "Good News and Bad News: Representation Theorems and Applications," *Bell Journal of Economics* **12**: 380-391
- [12] Mookherjee, D. and I. Png (1989), "Optimal Auditing, Insurance and Redistribution," *Quarterly Journal of Economics* **104**: 399-415.
- [13] Norris, R.S., H. M. Kristensen and J. Handler (2002), "Pakistan's Nuclear Forces, 2001," *Bulletin of the Atomic Scientists* **58**: 70-71
- [14] Norris, R.S. and H. M. Kristensen (2005a), "North Korea's Nuclear Program, 2005," *Bulletin of the Atomic Scientists* **61**: 64-67
- [15] Norris, R.S. and H. M. Kristensen (2005b), "India's Nuclear Forces, 2005," *Bulletin of the Atomic Scientists* **61**: 73-75
- [16] Schelling, T.C. (1960) *The Strategy of Conflict* (Cambridge, MA: Harvard University Press)
- [17] Schrage, M. (2003), "No Weapons, No Matter. We called Saddam's Bluff," Op-Ed in Washington Post, May 11, 2003.
- [18] Sobel, J. (1992), "How and (when) to Communicate with Enemies," in *Equilibrium and Dynamics*, edited by M. Majumdar, Macmillan, pages 307-321.

- [19] Takeyh, R. (2005), “WMD, Terrorism and Proliferation,” Testimony Before Subcommittee on Prevention of Biological and Nuclear Attack, Committee on Homeland Security, September 8, 2005
- [20] Thucydides (1972) *The History of the Peloponnesian War* (London: Penguin Classics)
- [21] Townsend, R. M. (1979), “Optimal contracts and costly state verification, *Journal of Economic Theory* **61**: 265-298.
- [22] *Treaty on the Non-proliferation of Nuclear Weapons*, <http://www.iaea.org/Publications/Documents/Treaties/npt.html>

## 6 Appendix

In this appendix, we characterize the set of all cheap-talk equilibria. Player A sends a message  $m \in M$ , where  $M$  is an arbitrary message space. Of course, the labelling of messages is arbitrary.

**Proposition 8** *All equilibria with effective cheap-talk can be replicated by using just two messages.*

**Proof.** Without loss of generality, we may assume player A’s types do not randomize over messages. Let  $A(m)$  be the set of types of player A who send message  $m$  in equilibrium. (That is,  $a \in A(m)$  if type  $a$  sends  $m$ .) Without loss of generality, we may assume  $A(m) \neq \emptyset$  for all  $m \in M$  (since a message which is never sent can be dropped). We know that the crazy type invests with probability 1. Let  $x(m)$  be the probability that the normal type of player B invests when player A sends  $m \in M$ . By Bayes’ rule, the probability that B is crazy conditional on being armed and message  $m$  having been sent is

$$\tau(m, \text{armed}) \equiv \frac{\tau}{\tau + (1 - \tau)x(m)}.$$

Notice that  $\tau(m, \text{armed}) \geq \tau$ . The set of messages that minimize  $x(m)$  is denoted  $M^C \subseteq M$ . Let  $M^T \equiv M \setminus M^C$ . By definition, if  $m^C \in M^C$  and  $m^T \in M^T$ , then  $x(m^C) < x(m^T) \leq 1$ . If communication is effective in influencing B’s investment decision, some types of A must send a message in  $M^C$ , and some types must send a message in  $M^T$ .

When B is unarmed, whether he is normal or crazy is payoff irrelevant, so there is no reason to distinguish the unarmed normal type from the unarmed crazy type. Abusing terminology, let B's *ex post type* be denoted  $t \in \{z, n, u\}$ , where  $n$  denotes that B is *armed and normal*,  $z$  that B is *armed and crazy*, and  $u$  denotes that B is *unarmed*. Let  $I(m, t, \varepsilon)$  be the probability that player B allows inspections following message  $m$ , when his ex post type is  $t \in \{z, n, u\}$  and the cost of inspection is  $\varepsilon \in [0, \bar{\varepsilon}]$ . This formulation is without loss of generality, because all unarmed players must make the same decision at the inspection stage (for the same  $\varepsilon$ ).

Conditional only on  $m$  and  $t$ , the probability of inspections is

$$I(m, t) \equiv \int_0^{\bar{\varepsilon}} I(m, t, \varepsilon) h(\varepsilon) d\varepsilon.$$

The probability that B is crazy conditional on being armed and allowing inspections following message  $m \in M$  is

$$\tau(m, \text{allow}, \text{armed}) \equiv \frac{\tau I(m, z)}{\tau I(m, z) + (1 - \tau) x(m) I(m, n)}.$$

The probability that B is crazy conditional on being armed and refusing inspections after message  $m \in M$  is

$$\tau(m, \text{refuse}, \text{armed}) \equiv \frac{\tau (1 - I(m, z))}{\tau (1 - I(m, z)) + (1 - \tau) x(m) (1 - I(m, n))}.$$

Also, let  $att(m, \text{allow}, \text{armed})$  be the probability that A attacks after message  $m$  if B allows inspections which reveal he is armed. Similarly,  $att(m, \text{allow}, \text{unarmed})$  is the probability that A attacks if B allows inspections which reveal he is unarmed, and  $att(m, \text{refuse})$  is the probability that A attacks if B refuses inspections. The proof has 14 steps.

*Step 1:* Suppose after message  $m$ , there is positive probability that inspections reveal B is armed. Then,

- (a)  $att(m, \text{allow}, \text{armed}) < att(m, \text{refuse})$ ;
- (b) if type  $a \in A(m)$  sends message  $m$  and then attacks when inspections reveal B is armed, then he must also attack if inspections are refused; and
- (c)  $\tau(m, \text{allow}, \text{armed}) \geq \tau(m, \text{armed}) \geq \tau(m, \text{refuse}, \text{armed})$  (the inequalities are strict if B refuses inspections with positive probability when armed).

*Proof:* (a) Both armed types of B want to minimize the probability of attack as  $\delta_t \geq 0 > -\alpha + \gamma$ . Also, inspections are costly for B. Hence, when

armed, B will allow inspections after message  $m \in M$  only if they lower the probability of attack.

(b) If  $a \in A(m)$  and type  $a$  attacks if inspections are refused, then if  $a' > a$  and  $a' \in A(m)$ , type  $a'$  also attacks. Hence, for any message  $m \in M$ , there is a cut-off type  $\tilde{a}$  such that if inspections are refused, then type  $a \in A(m)$  attacks if and only if  $a \geq \tilde{a}$ . Similarly, there is a type  $\hat{a}$  such that if inspections reveal that B is armed, then type  $a \in A(m)$  attacks if and only if  $a \geq \hat{a}$ . Part (a) implies  $\hat{a} > \tilde{a}$ , which implies (b).

(c) The armed type  $t \in \{n, z\}$  with cost shock  $\varepsilon$  is willing to allow an inspection if and only if

$$\begin{aligned} & att(m, \text{allow}, \text{armed}) (-\alpha + \gamma) + (1 - att(m, \text{allow}, \text{armed})) \delta_t - \varepsilon \\ & \geq att(m, \text{refuse}) (-\alpha + \gamma) + (1 - att(m, \text{refuse})) \delta_t. \end{aligned}$$

If this condition holds for  $t = n$ , then it holds strictly for  $t = z$ , because  $\delta_z > \delta_n$  and  $1 - att(m, \text{allow}, \text{armed}) > 1 - att(m, \text{refuse})$  by part (a). As crazy types are more willing to allow inspections than normal types, (c) is proved. This completes the proof of step 1.

The remaining steps will establish that, without loss of generality, we can assume  $M^C = \{m^C\}$  and  $M^T = \{m^T\}$ . Moreover, we will show that

$$\begin{aligned} I(m^T, z) &= I(m^T, n) = 1, \\ I(m^T, u) &= I(m^C, u) = 0 \end{aligned}$$

and either

$$I(m^C, z) = I(m^C, n) = 0$$

or

$$1 = I(m^C, z) > I(m^C, n) > 0.$$

This corresponds with the communication equilibrium from Section 4, if  $m^C$  is interpreted as the conciliatory message and  $m^T$  as the tough message. It is then easy to show that investment and attack decisions must also be the same.

*Step 2:* Doves and hawks only send messages in  $M^C$ .

*Proof:* Suppose  $a \in A(m)$ , where  $a$  is a dove or a hawk. We claim  $m \in M^C$ . Notice that if, along the equilibrium path following message  $m$ , type  $a$  either *never* attacks or *always* attacks, then clearly  $m \in M^C$ , for otherwise

type  $a$  could increase his expected payoffs by reducing the probability that B invests. There are three cases to consider.

*Case 1:* Suppose inspections never occur on the equilibrium path after message  $m$  is sent.

For a dove with  $a \leq 0$ , the expected payoff from not attacking is

$$-\sigma(\tau d_z + (1 - \tau)x(m)d_n) \geq -\sigma(\tau d_z + (1 - \tau)d_n) > a - \sigma c$$

by Assumption 1, so his best-response is to never attack. Hence,  $m \in M^C$ .

For a hawk with  $a \geq c - (\tau d_z + (1 - \tau)d_n)$ , the expected payoff from attacking is

$$\begin{aligned} & a - \sigma(\tau + (1 - \tau)x(m))c \\ = & (1 - \sigma(\tau + (1 - \tau)x(m)))a + \sigma(\tau + (1 - \tau)x(m))(a - c) \geq \\ & (1 - \sigma(\tau + (1 - \tau)x(m)))a - \sigma(\tau + (1 - \tau)x(m))(\tau(m, \text{armed})d_z + (1 - \tau(m, \text{armed}))d_n) \\ = & (1 - \sigma(\tau + (1 - \tau)x(m)))a - \sigma(\tau d_z + (1 - \tau)x(m)d_n) \\ > & -\sigma(\tau d_z + (1 - \tau)x(m)d_n) \end{aligned}$$

as  $a > 0$  and  $\tau(m, \text{armed}) \geq \tau$ . Hence, a hawk's best-response is to always attack. Hence,  $m \in M^C$ .

*Case 2:* Suppose inspections always occur on the equilibrium path after message  $m$  is sent.

Suppose type  $a \in A(m)$  is a dove. If inspections reveal that B is unarmed, the dove does not attack. If in addition he doesn't attack when inspections reveal that B is armed, then he never attacks, so  $m \in M^C$ . Suppose instead that type  $a$  attacks if inspections reveal B is armed. Then all types  $a' \in A(m)$  with  $a' \geq a$  also attack if inspections reveal B is armed. By part (a) of step 1, an armed type must expect that an attack occurs with higher probability if he refuses inspections. Hence, there must be some types  $a' \in A(m)$  with  $a' < a$  who do not attack when inspections reveal B is armed. Hence, there is  $a' \in A(m)$  who never attacks, so  $m \in M^C$ .

Suppose type  $a \in A(m)$  is a hawk. Following message  $m$ , type  $a$  attacks inspections reveal B is unarmed. If inspections reveal B is armed, the payoff from attacking is

$$a - c \geq -(\tau d_z + (1 - \tau)d_n) \geq -(\tau(m, \text{armed})d_z + (1 - \tau(m, \text{armed}))d_n)$$

and hence  $a$  always attacks. Therefore,  $m \in M^C$ .

*Case 3:* Suppose inspections sometimes (but not always) occur on the equilibrium path after message  $m$  is sent.

*Sub-case 3.1:* If B always allows inspections when armed, then B's equilibrium strategy fully reveals his capabilities with probability one. Hence, the argument in Case 2 shows  $m \in M^C$ .

*Sub-case 3.2:* Suppose armed types sometimes accept and sometimes reject inspections on the equilibrium path in response to message  $m \in M$ .

Suppose type  $a \in A(m)$  is a dove. Then type  $a$  does not attack if inspections reveal B is unarmed. Suppose type  $a$  attacks if inspections are refused. This implies that

$$-\lambda(\tau(m, \text{refuse, armed})d_z + (1 - \tau(m, \text{refuse, armed}))d_n) \leq a - \lambda c < -\lambda c \quad (25)$$

where  $\lambda$  is the probability that B is armed conditional on refusing inspections. (The first inequality of (25) says that type  $a$  prefers to attack, the second is due to  $a < 0$ .) Notice that (25) implies that

$$-c > -(\tau(m, \text{refuse, armed})d_z + (1 - \tau(m, \text{refuse, armed}))d_n). \quad (26)$$

We then have

$$a - c \geq -(\tau(m, \text{refuse, armed})d_z + (1 - \tau(m, \text{refuse, armed}))d_n) > -(\tau(m, \text{allow, armed})d_z + (1 - \tau(m, \text{allow, armed}))d_n) \quad (27)$$

The first inequality in (27) is due to the fact that, since the first inequality in (25) holds for  $\lambda \leq 1$ , (26) implies it also holds when  $\lambda$  is replaced by 1. The second inequality in (27) is due to  $\tau(m, \text{allow, armed}) > \tau(m, \text{refuse, armed})$ , from part (c) of step 1. Now, (27) implies that type  $a$  strictly prefers to attack when B allows inspections which reveal that he is armed. Hence, any dove in  $A(m)$  who attacks when inspections are refused also attacks when inspections show B is armed. More generally, if type  $a \in A(m)$  attacks when inspections show B is armed, then any type  $a' > a$  in  $A(m)$  also attacks when inspections show B is armed. But this implies  $\text{att}(m, \text{allow, armed}) \geq \text{att}(m, \text{refuse})$ , contradicting part (a) of step 1. This contradiction proves that type  $a$  does not attack if inspections are refused. If type  $a$  attacks when inspections reveal that B is armed then, by step 1 (b), he must also attack if inspections are refused, which contradicts the previous sentence. So in fact a dove  $a \in A(m)$  never attacks, hence  $m \in M^C$ .

Suppose type  $a \in A(m)$  is a hawk. Then,

$$a - c \geq -(\tau d_z + (1 - \tau) d_n) > -(\tau(m, \text{allow, armed})d_z + (1 - \tau(m, \text{allow, armed}))d_n) \quad (28)$$

(The first inequality is due to the definition of hawk, the second is due to  $\tau(m, \text{allow}, \text{armed}) > \tau(m, \text{armed}) \geq \tau$  from step 1 (c).) Now (28) implies that type  $a$  attacks if inspections reveal B is armed. Part (b) of step 1 implies type  $a$  also attacks if inspections are refused. Since hawks certainly attack if inspections reveal B is unarmed, type  $a$  always attacks. Hence  $m \in M^C$ .

*Sub-case 3.3:* Finally, suppose B never allows inspections when armed, after  $m$  was sent. Then he must sometimes allow inspections when unarmed, otherwise we are in case 1. Hence, the probability of attack must be lower after inspections reveal B is unarmed than after he refuses inspections:

$$\text{att}(m, \text{allow}, \text{unarmed}) < \text{att}(m, \text{refuse}). \quad (29)$$

By an argument like in the proof of step 1 part (b), (29) implies that if there is any type in  $A(m)$  who attacks when inspections reveal B is unarmed, he must also attack when inspections are refused, hence he always attacks. In this case,  $m \in M^C$ . Suppose instead that no type in  $A(m)$  attacks when inspections reveal B is unarmed, i.e., only doves send message  $m$ . But then in response to message  $m$ , a normal type of B should refrain from investing, and then reveal that he is unarmed. Therefore,  $x(m) = 0$ , which certainly implies  $m \in M^C$ .

*Step 3:*  $I(m^T, u) = 0$  for all  $m^T \in M^T$ .

*Proof:* By step 2, message  $m^T \in M^T$  reveals that A must be an opportunist. Therefore, if B is unarmed, then he will refuse inspections following message  $m^T \in M^T$ . (If he were to allow inspections, he would pay  $\varepsilon > 0$  only to be attacked for sure.)

*Step 4:*  $I(m^T, n) = I(m^T, z) = x(m^T) = 1$  for all  $m^T \in M^T$ . If message  $m^T \in M^T$  is sent and there is no inspection, then an attack occurs with probability one.

*Proof:* Inspections must occur with positive probability following  $m^T$ , or else all types of A would prefer to send  $m^C \in M^C$ . Since  $I(m^T, u) = 0$  by step 3, whenever inspections occur they must reveal that B is armed, so step 1 applies. Specifically, if some type  $a \in A(m^T)$  attacks after the inspections, then he must also attack if inspections are refused, by step 1 (b). That means he always attacks (since  $I(m^T, u) = 0$ ), so he is better off sending  $m^C \in M^C$ , a contradiction. Hence, no type in  $A(m^T)$  attacks after inspections have revealed B is armed. If some type  $a \in A(m^T)$  doesn't attack when inspections are refused, then this type never attacks (since  $I(m^T, u) = 0$ ), so he is better off sending  $m^C \in M^C$ , a contradiction. Hence, when inspections

are refused, the probability of attack is one. On the other hand, as we have seen, if inspections reveal B is armed then the probability of attack is zero. This clearly means B strictly prefers to allow inspections when armed, hence  $I(m^T, n) = I(m^T, z) = 1$  for any  $m^T \in M^T$ . It also implies B strictly prefers to invest when he hears message  $m^T \in M^T$ , so  $x(m^T) = 1$ . Indeed, if after hearing message  $m^T \in M^T$  player B invests and allows inspections if and only if he is armed, his payoff is  $\sigma(\delta_t - E\varepsilon) - (1 - \sigma)\alpha - k$ . By not investing he gets  $-\alpha$ . Assumption 3 implies that he prefers to invest.

*Step 5:* We can assume, without loss of generality, that  $M^T = \{m^T\}$  is a singleton.

*Proof:* Steps 2 and 4 imply that any type of player A who sends a message  $m^T \in M^T$  must be an opportunist who attacks if and only if inspections reveal no arms or inspections are refused. Also,  $x(m^T) = 1$  for all  $m^T \in M^T$ . Therefore, all messages in  $M^T$  lead to the same outcome, so we may as well assume there is only one such message.

*Step 6:* Following any  $m^C \in M^C$ , there is a positive probability that B is armed and refuses inspection.

*Proof:* Since the only reason to send the message  $m^T \in M^T$  is to improve the chance of an inspection, steps 3 and 4 imply we cannot have  $I(m^C, n) = I(m^C, z) = 1$ .

*Step 7:* Doves never attack on the equilibrium path and hawks always attack.

*Proof:* Consider the doves. By step 2, we know they send only messages in  $M^C$ . In the proof of step 2, case 1, we established that if inspections never occur on the equilibrium path, doves do not attack. So assume inspections occur with positive probability on the equilibrium path following  $m^C$ . First, suppose  $I(m^C, u) > 0 = I(m^C, z) = I(m^C, n)$ . Now, if inspections reveal B is unarmed, then all types with  $a > 0$  attack. For the unarmed to be willing to allow inspections, the inspections must strictly reduce the probability of attack. This implies that some doves ( $a < 0$ ) who send  $m^C$  must attack if inspections are refused. These doves get no more than  $-\sigma(\tau + (1 - \tau)x(m^C))c$  from sending  $m^C$ , since they end up attacking whenever B is armed. If instead they send  $m^T$  and never attack, they get  $-\sigma(\tau d_z + (1 - \tau)d_n)$ . Thus, for them to prefer to send  $m^C$ , we need

$$-\sigma(\tau + (1 - \tau)x(m^C))c \geq -\sigma(\tau d_z + (1 - \tau)d_n) \quad (30)$$

If type  $a > 0$  sends  $m^T$ , as shown above, in equilibrium he will attack whenever B is unarmed, and he gets

$$(1 - \sigma)a - \sigma(\tau d_z + (1 - \tau)d_n) < a - \sigma(\tau + (1 - \tau)x(m^C))c$$

The inequality uses  $1 - \sigma < 1$  and (30). But, the right hand side is what he gets if he sends  $m^C$  and always attacks. Thus, no type  $a > 0$  will send  $m^T$ , so  $m^T$  is never sent in equilibrium, a contradiction. Therefore,  $I(m^C, u) > 0 = I(m^C, z) = I(m^C, n)$  is impossible. Thus, if inspections sometimes occur, they must sometimes occur when B is armed, so we can use step 1. By step 6, they must also be refused with positive probability when B is armed. Suppose type  $a = 0$  sends  $m^C$  and inspections are refused. Let  $\tilde{\sigma}$  denote the probability B is armed, conditional on inspections being refused after message  $m^C \in M^C$ . Then, after the refusal, type  $a = 0$  prefers to attack if and only if

$$-\tilde{\sigma}c \geq -\tilde{\sigma}(\tau(m, \text{refuse})d_z + (1 - \tau(m, \text{refuse}))d_n) \quad (31)$$

We will show this inequality leads to a contradiction. Step 1 (c) implies  $\tau(m, \text{refuse}) \leq \tau(m, \text{allow}, \text{armed})$ . Therefore, (31) implies

$$-c \geq -(\tau(m, \text{allow}, \text{armed})d_z + (1 - \tau(m, \text{allow}, \text{armed}))d_n)$$

so (31) implies type  $a = 0$  also prefers to attack if inspections reveal B is armed. Hence, type  $a = 0$  prefers to always attack, which is to say

$$-\sigma(\tau + (1 - \tau)x(m^C))c \geq -\sigma(\tau d_z + (1 - \tau)d_n).$$

But then, as  $1 - \sigma < 1$ , for opportunists with  $a > 0$  we have

$$a - \sigma(\tau + (1 - \tau)x(m^C))c > (1 - \sigma)a - \sigma(\tau d_z + (1 - \tau)d_n). \quad (32)$$

The right-hand-side of (32) is the expected payoff to an opportunist from sending message  $m^T \in M^T$ . Hence, (32) implies no opportunist sends a message in  $M^T$  and in fact all types send messages in  $M^C$ . But this contradicts our assumption that the equilibrium displays effective communication. This contradiction shows (31) cannot hold, so type  $a = 0$  strictly prefers *not* to attack if inspections are refused, and hence so do all types with  $a < 0$ . If there is positive probability that inspections reveal B is armed, then if type  $a < 0$  attacks in this case, he must also attack if inspections are refused, by

step 1 (b), which contradicts the previous sentence. Thus, if type  $a < 0$  sends  $m^T$ , he will not attack if inspections are refused, or if they reveal that B is armed. Since he will certainly not attack if inspections show B is unarmed, type  $a < 0$  will never attack. A similar argument shows that hawks always attack on the equilibrium path.

*Step 8:*  $I(m^C, u) = 0$  for all  $m^C \in M^C$ .

*Proof:* If message  $m^C \in M^C$  is sent and inspections are refused, hawks attack but not doves (by step 7). Clearly, if B hears message  $m^C \in M^C$  and is unarmed, he has no reason to allow inspections. Indeed, inspections cannot convince the hawks not to attack, and if some opportunist also sends  $m^C$ , he will attack if inspections reveal that B is unarmed. Therefore,  $I(m^C, u) = 0$ .

*Step 9:*  $0 < x(m^C) < 1$  for  $m^C \in M^C$ , so B's normal type is indifferent between investing and not investing when he hears  $m^C$ .

*Proof:* Suppose  $x(m^C) = 0$  for  $m^C \in M^C$ . Then, after message  $m^C$ , only crazy types are armed. Clearly, any type  $a > 0$  prefers to send  $m^C$  rather than  $m^T$ , in view of  $m^T = 1$ . But then, by step 2, no type sends  $m^T$ , which contradicts communication being effective. Thus,  $x(m^C) = 0$  is impossible. On the other hand,  $x(m^C) < x(m^T) = 1$ .

*Step 10:* There are cut-off points  $a'$  and  $a''$ , where  $0 < a' \leq a'' < c - \tau d_z - (1 - \tau) d_n$ , such that A sends  $m^T$  if  $a \in (a', a'')$ . He sends some  $m^C \in M^C$  if  $a < a'$  or  $a > a''$ .

*Proof:* The proof of Step 7 established that type  $a = 0$  strictly prefers not to attack if inspections are refused on the equilibrium path. Some “weak” opportunists with  $a$  close to zero must have the same strict preference as type  $a = 0$ . Let  $a' > 0$  be the supremum of all such types. All types such that  $a < a'$  must send  $m^C \in M^C$  and then not attack if there are no inspections. Similarly, it can be shown that “tough” opportunists with types just less than  $c - \tau d_z - (1 - \tau) d_n$  must be sending messages in  $M^C$  and attacking if player B refuses inspections... Let  $a''$  be the infimum of all such types. All types such that  $a > a''$  must send  $m^C$  and then attack if there are no inspections. Necessarily,  $0 < a' \leq a'' < c - \tau d_z - (1 - \tau) d_n$ .

*Step 11:* Either  $I(m^C, n) = I(m^C, z) = 0$  for all  $m^C \in M^C$ , or there is  $I^* \in (0, 1)$  such that for all  $m^C \in M^C$ ,  $I(m^C, n) = I^*$  and  $I(m^C, z) = 1$ .

*Proof:* Suppose there is  $m^C$  such that  $I(m^C, t) > 0$  for some armed  $t \in \{n, z\}$ . We cannot have  $I(m^C, z) = I(m^C, n) = 1$ , by step 6. By an argument as in step 1 (c), crazy types are more willing to allow inspections

than normal types, so we must have  $I(m^C, z) > I(m^C, n)$ . Now, if type  $z$  is willing to allow inspections, they must reduce the probability of attack. Hawks always attack by step 7. Therefore, if inspections reveal that B is armed, some “tough” opportunists ( $a'' < a < c - \tau d_z - (1 - \tau) d_n$ ) must not attack. If  $I(m^C, n) = 0$ , then when inspections reveal B is armed A can infer that B is crazy, so all opportunists would attack, contradicting the previous sentence. Therefore,  $0 < I(m^C, n) < I(m^C, z) = 1$ . By definition,  $x(m^C)$  is the same for all  $m^C \in M^C$ . This implies we must have  $I(m^C, n)$  constant for all  $m^C$  that are sent in equilibrium...

*Step 12:* We can assume, without loss of generality, that  $M^C = \{m^C\}$  is a singleton.

*Proof:* Step all messages in  $M^C$  yields the same outcome, so we may assume there is only one such message.

*Step 13:* If  $I(m^C, n) = I(m^C, z) = 0$  for all  $m^C \in M^C$  then the cut-off types  $a'$  and  $a''$  (defined in step 9) and the investment probability  $x = x(m^C)$ , are determined by equations (19), (20) and  $\Psi(x) = 0$  (where  $\Psi$  is defined by (23)).

*Proof:* It must be the case if  $a' < a''$ ,  $a'$  and  $a''$  are indifferent between reporting  $m^C$  and  $m^T$ . These indifference conditions yield (19) and (20). By step 8, B must be indifferent between investing and not investing when he hears  $m^C$ . The equation  $\Psi(x) = 0$  is this indifference condition.

Step 13 proves that any equilibrium with  $I(m^C, n) = I(m^C, z) = 0$  for all  $m^C \in M^C$  is outcome-equivalent to the communication equilibrium described in Section 4. ■