



**CMS-EMS**  
**Center for Mathematical Studies in Economics  
And Management Science**

Discussion Paper #1572

**Achieving Cooperation  
Under Privacy Concerns**

Wioletta Dziuda  
Ronen Gradwohl  
Northwestern University

First version: April 2012  
This version: November 5, 2013

*Keywords:* Cooperation, Privacy, Communication

JEL classification: D80



**NORTHWESTERN  
UNIVERSITY**



# Achieving Cooperation under Privacy Concerns\*

Wioletta Dziuda<sup>†</sup>      Ronen Gradwohl<sup>‡</sup>

First version: April 2012

This version: November 5, 2013

## Abstract

Two players choose whether to cooperate on a project. Each of them is endowed with some evidence, and if both possess a sufficient amount then cooperation is profitable. In order to facilitate cooperation the players reveal evidence to one another. However, some players are concerned about privacy, and so revelation of evidence that does not result in cooperation is costly.

We show that in equilibrium evidence can be exchanged both incrementally and all at once, and identify conditions under which the different rates of evidence exchange are optimal.

---

\*We are grateful to the participants of seminars at Northwestern, UCSD, and Penn State, as well as the Midwest Economic Theory Meetings in St. Louis.

<sup>†</sup>Kellogg School of Management, Northwestern University, Evanston, IL 60208, USA. E-mail: [wdziuda@kellogg.northwestern.edu](mailto:wdziuda@kellogg.northwestern.edu).

<sup>‡</sup>Kellogg School of Management, Northwestern University, Evanston, IL 60208, USA. E-mail: [r-gradwohl@kellogg.northwestern.edu](mailto:r-gradwohl@kellogg.northwestern.edu).

# 1 Introduction

When two parties communicate in an attempt to undertake a joint venture, the conventions and protocols that structure their communication may be formed by a variety of factors. In this paper we analyze the interplay of two such factors: On the one hand, in order to cooperate successfully a party must communicate some proprietary information that is necessary for the venture. On the other hand, parties may have privacy concerns: If the joint venture fails to materialize, a party may be adversely affected by the other's use of the revealed information.

Consider the following examples of communication with such privacy concerns: Two firms with complementary expertise wish to cooperate on a project. The execution and success of the project depend on firms sharing their expertise and ideas. However, some firms' level of expertise and novelty of ideas may be too poor to permit successful cooperation. If this fact becomes clear in the process of communication, the project is abandoned. In such an event, however, a firm that revealed promising ideas prior to the abandonment may regret doing so. Anticipating this, firms may want to structure their communication in ways that minimize the harm sustained in case cooperation fails.

Next, consider two shady characters who wish to engage in a less-than-legal venture. They exchange plans for a potential criminal scheme, references to criminal connections that may be useful in their venture, and descriptions of other activities to which only the criminal underworld is privy. However, some characters may be undercover cops—they are uninterested in a criminal venture, but rather in obtaining information from the criminal that may lead to an arrest. So while the potential profits from the venture may render information exchange appealing, the characters' concerns for privacy drive them to structure communication in a way that minimizes the amount of incriminating information revealed to undercover agents.

As a third example, consider two researchers, a theoretician and an empiricist, who are interested in joining forces on a project. They engage in communication

to convey the content of their research and write a joint paper. At the same time, they face uncertainty about the viability of their potential research partner. Does the theoretician actually have a sound and reasonable model, and does she have theorems with correct proofs that fit the project's aims? Does the empiricist have sufficient data, and do those data support the project's aims? If the answer is no, the researcher may be unwilling to reveal her own ideas, as those might potentially be exploited by the other.

In this paper, we analyze the tradeoff between the two conflicting forces demonstrated in the examples—the necessity of information exchange versus the concern for privacy—and examine its effect on the structure of communication. In particular, what is the optimal rate of information exchange? Is optimal communication incremental, with parties revealing little bits of information in alternating fashion? Or is optimal communication simple, with one party revealing all her information at once?

Our main result is that both modes of communication—incremental and simple—may be optimal in equilibrium. Which of the two obtains depends on the order in which information must be revealed. Suppose that the initial pieces of information to be exchanged are unlikely to reveal the viability of the opponent or that they tend to inflict a relatively high cost due to privacy loss. In this case, the optimal mode of communication is simple, which one player revealing all evidence first. By a similar intuition, incremental information exchange is optimal when the initial pieces of information to be exchanged are relatively likely to reveal the viability of the opponent or when they inflict little harm due to privacy loss.

In this paper, we take the order of information exchange as given. Such an assumption is plausible in many applications in which information must be revealed in a predetermined order. For example, a theoretician describing a proof must reveal lemmas that build one on top of another. However, in other applications parties may have more flexibility as to the order in which information is

presented: a firm can reveal its financial statements to potential partners before or after allowing them to visit its factories, and criminals may agree on what proofs of viability to present first. For such cases, our analysis sheds light on the optimal order in which information should be revealed. We show that privacy leakage is smaller when incremental exchange is optimal than when simple exchange is optimal. Hence, when they have the flexibility to do so, parties should order information in such a way that optimal communication is incremental: information that is less valuable or more likely to demonstrate viability should be revealed first.

**Organization** Immediately following is a brief survey of the related literature. Sections 2 and 3 contain the model and its analysis, the latter of which includes our main results about the optimality of incremental and simple communication. In Section 4 we then discuss the robustness of the results to our assumptions. Finally, most proofs are deferred to the Appendix.

## 1.1 Related Literature

This paper is part of a large literature on strategic information exchange that was pioneered by Crawford and Sobel (1982) (cheap-talk), Milgrom (1981), and Milgrom and Roberts (1986) (verifiable information). Our notion of information is modeled after Shin (1994) and Dziuda (2011), in that players must reveal their information truthfully, but can obstruct their type by withholding some information. Unlike most of the communication literature (exceptions include Li et al. (2001)), in our paper both players have private information about the relevant state variable (their type) and both communicate this information. Moreover, in our paper the messages sent affect players' utilities. This is similar to Kartik (2009), who assumes that lying is costly, with the difference that in our paper revealing (by assumption truthfully) information is costly.

The paper that is closest to ours is Augenblick and Bodoh-Creed (2013). In

that paper, developed concurrently with and independently of ours, each party has a privately observed type and wishes to find a matching partner, but prefers to confuse non-matching partners about her type. The authors assume that pieces of information are heterogenous and focus on the order in which they should be revealed. We analyze a somewhat orthogonal problem by taking the order of information exchange as given (perhaps because it is agreed upon beforehand) and focusing on whether information exchange should be simple or incremental. This question is minor in Augenblick and Bodoh-Creed (2013), because in their setting communication is essentially one-sided: One player reveals information about her type, and the other confirms whether her type matches or not. Since Augenblick and Bodoh-Creed (2013) search for the sender-optimal equilibrium, incremental communication is optimal: at the first point at which the receiver does not confirm the match, communication stops. Hence, incremental exchange allows the sender to avoid revealing her entire type. In our model communication must be two-sided as each of the viable types possesses different information. As a result, the question of whether incremental communication is better arises, and the answer turns out to depend on the nature of the evidence. While our results do shed some light on the optimal order of information exchange, these implications are somewhat incomparable with the results of Augenblick and Bodoh-Creed (2013) because the modeling of privacy concerns is very different in the two papers.

Hörner and Skrzypacz (2011) analyze a problem in which an uninformed principal wishes to acquire information from a possibly informed agent. The principal cares about money (which translates to a form of privacy concerns), but the agent does not have privacy concerns. Hörner and Skrzypacz (2011) show that in the equilibrium that maximizes the surplus of the principal and the informed agent, information is revealed and payments are made gradually.

Information in our model can be interpreted as money or effort, and hence communication can be viewed as contributions to a common project. In most models on this topic, there is a free-riding problem: Each player has an incentive

to withhold her contributions and wait for the opponent to finance the public good. This literature finds that gradualism—splitting contributions into smaller pieces and contributing in an alternating fashion—may ameliorate this problem (see Lockwood and Thomas (2002), Marx and Matthews (2000), Compte and Jehiel (2004) in the context of public good, and Pitchford and Snyder (2004) in the context of the hold-up problem).<sup>1</sup> In contrast to this literature, there is no free-riding in our model as players' actions are not substitutes but complements. Under symmetric information, if both players are able to contribute to the project, they would do so. Hence, gradualism in our model stems from the uncertainty about the ability of the opponent to contribute.

Watson (1999, 2002) obtains that gradualism may be optimal in partnerships if asymmetric information is present. In these papers, two players are involved in a partnership: High types want to stay in the partnership forever, while low types have an incentive to exit unless the partnership level increases quickly. Watson (1999, 2002) shows that high types prefer to start at a low level of partnership and increase it slowly, as this encourages the low types to exit early, when the stakes are low. In our paper, the incremental exchange also allows the low (unviable) types to drop out (being discovered in our case) at the moment at which the opponent invested only little (revealed little information in our case). However, unlike in Watson (1999, 2002), players are allowed to participate in the exchange in an asymmetric fashion. Hence, incremental information exchange corresponds to both players investing little in a partnership, while simple information exchange corresponds to one player investing more than the opponent. As a result, we obtain that gradualism may not be optimal.

Also related is the literature on sustaining conversations, including Stein

---

<sup>1</sup>In Admati and Perry (1991) gradualism in contributions comes from the convexity of the cost function and would be optimal even with a single player. Moreover, as Compte and Jehiel (2003) argue, the insight of Admati and Perry (1991) about gradualism is sensitive to the symmetry assumption. In contrast, in our model gradualism arises solely as result of strategic interactions and would hold also if players were asymmetric.

(2008) and Ganglmair and Tarantino (2013). While these papers also have an element of privacy concerns, the driving force behind their models is that conversation generates new ideas, and hence players converse to further develop new information.

Our paper is somewhat related to the computer science literature on exchange protocols, which also derived the insight that incremental communication can be beneficial. In particular, Blum (1983), Damgård (1995), and Bardsley et al. (2008) show that incremental communication can facilitate the exchange of secrets. However, a critical element of these papers is that players can, for some cost, discover the opponent's secret even in the absence of communication. This element, which is absent from our paper and others in the strategic communication literature, is crucial for the equilibrium analysis of the aforementioned exchange protocols.

Finally, our paper is related to the cryptography literature, and in particular to zero-knowledge proofs (pioneered by Goldwasser et al. (1989)) and secure 2-party computation (introduced by Yao (1982)). The tools developed in this field allow computationally-bounded players to jointly compute a function of their respective private information, while maintaining the privacy of this information. However, in many economic applications information cannot be generically encoded and must be seen by the opponent to be verified, which would render these tools inapplicable. Furthermore, in some applications it is simply not feasible or too costly to run a cryptographic protocol to reach agreement.

## 2 The Model

**Players and their types** There are two players  $\{1, 2\}$ , which we typically denote by  $i \in \{1, 2\}$  and  $j \stackrel{\text{def}}{=} 3 - i$ , each of which has a type  $\tau_i$ . Each player possesses a unit of evidence. This evidence can potentially lead to a successful project, in which case we call the player *viable* ( $\tau_i = V$ ), or not, in which case we

call the player *unviable* ( $\tau_i = U$ ). The type of a player is her private information. The prior probability that a player is viable is  $p$  and is independent across players.

Evidence is to be interpreted as a code, a recipe, or a proof that takes a fixed amount of time (or space) to transmit, but can be divided into smaller pieces, each of which can be transmitted in correspondingly less time (or space).

**Game** There are possibly infinitely many rounds of communication. In each round, one player is called upon to speak, and this happens in an alternating fashion with player 1 moving first. In each round  $t$ , the speaking player  $i$  chooses the amount of new evidence to disclose in that round. Formally, let  $N_t^i$  be the amount of evidence disclosed by player  $i$  up to round  $t$  (including  $t$ ). In each round  $t$ , the speaking player  $i$  chooses  $N_t^i \in [N_{t-2}^i, 1]$ , where  $N_{-1}^i = N_0^i \stackrel{\text{def}}{=} 0$ . That is, we are assuming that a player can withhold evidence, but cannot withdraw evidence already disclosed.

We want the model to capture the idea that a viable type needs to reveal the entire proof, code, or recipe in order to prove its viability. Unviable types are those whose proofs or recipes are incomplete or contain a fatal flaw. Hence, after an unviable player reveals a sufficient amount of the evidence, her type becomes known to the opponent. To this end, we assume that the unviable type of player  $i$  is characterized by a number  $K^i \in (0, 1)$ , which is her private information. If in some round  $t$ , the unviable type with  $K^i$  reveals  $N_t^i > K^i$ , player  $j$  receives a signal that the opponent is unviable:  $s_t^i = U$ . We assume that  $K^i$  is distributed with a strictly increasing, continuous distribution function  $F(K^i)$  and is independent of  $K^j$ . Note that we are implicitly assuming that an unviable player cannot fabricate the evidence of a viable one.<sup>2</sup>

---

<sup>2</sup>Note that  $K^i$  can be also interpreted as the highest amount of evidence that the unviable type possesses; for example, the number of lemmas proved. In this case after revealing  $K^i$  amount of evidence, player  $i$  would be unable to reveal any further evidence. We find it convenient to assume that  $i$  can continue revealing evidence, but since this evidence does not enter the payoffs, our assumption is without loss of generality.

A history of play up to and including round  $t$  is  $H_t = \{\{N_1^1, N_2^2, \dots, N_t^i\}, s_{t-1}^j, s_t^i\}$ , where  $s_t^i \in \{\emptyset, U\}$  and  $s_t^i = U$  means that player  $j$  received a signal that  $i$  is unviable and  $s_t^i = \emptyset$  means that player  $j$  did not receive this signal. A pure strategy of player  $i$  is a function that for each  $t$  in which  $i$  speaks maps  $H_{t-1}$  into  $N_t^i \in [N_{t-2}^i, 1]$ .<sup>3</sup>

We allow the players to split evidence into arbitrarily small pieces, but we want the number of such pieces to be finite. Hence, we place the following assumption on the set of strategies available to the players:

**Assumption 2.1** *A strategy  $\sigma_i$  of player  $i$  is admissible if for each strategy  $\sigma_j$  of player  $j$  there exists some  $T(\sigma_j)$  such that the following is satisfied: if  $i = 1$  ( $i = 2$ ), then for every odd (even)  $t > T(\sigma_j)$ , the history  $H_t$  generated by strategies  $\sigma_i$  and  $\sigma_j$  has the property that  $N_t^i = N_{t-2}^i$ .*

This assumption says that no matter how the game unravels, player  $i$  will stop revealing new evidence after a finite number of rounds.

For fixed strategies of the players, denote by  $N^i = \max_t N_t^i$  the largest amount of evidence revealed in the game by  $i$ . By Assumption 2.1,  $N^1$  and  $N^2$  exist.

**Payoffs** There exists a project that pays  $v > 0$  to each player if and only if both players are viable and share all their respective evidence: that is, if  $\tau_1 = \tau_2 = V$  and  $N^1 = N^2 = 1$ . In this case, we say that players *cooperate* on the project.

In addition to the payoff from the project, players obtain payoffs from the evidence exchanged in the game. We assume that only the evidence provided by the viable types is valuable, and discuss the robustness of our results to this assumption in Section 4. A viable player  $i$  who reveals  $N^i$  suffers a disutility

---

<sup>3</sup>It is standard in game theory to present dynamic games as evolving in physical time, and hence to describe the set of actions available to each player at each instance of time. Since time is continuous in our model, we find it more convenient to treat the uninterrupted flow of evidence as one round, use then rounds as different stages of the game, and describe the set of actions available in each round.

$h(N^i)$ , and a player  $i$  who receives  $N^j$  from a viable opponent benefits  $g(N^j)$ . Formally, the utilities of the viable and unviable types, respectively, are:

$$u_i^V = v \mathbf{1}_{(\tau_j=V, N^1=N^2=1)} + g(N^j) \mathbf{1}_{(\tau_j=V)} - h(N^i), \quad (1)$$

and

$$u_i^U = g(N^j) \mathbf{1}_{(\tau_j=V)}, \quad (2)$$

where the symbol  $\mathbf{1}_{(\cdot)}$  denotes the indicator function. Both  $g$  and  $h$  are strictly increasing, continuous, and take value 0 at 0. We normalize  $h(1) = 1$ .

**Equilibrium** The solution concept is a pure-strategy weak Perfect Bayesian Equilibrium (PBE).

Throughout the paper, we will primarily be interested in optimal modes of communication; in particular, in whether it is optimal for each player to reveal all evidence at once, or whether splitting evidence into finer pieces and revealing them in an alternating fashion can be welfare improving. Our optimality criterion will be the joint payoff of the viable types. If our game is part of a larger game in which players make costly investments to become viable, then maximizing the payoff of the viable types can be consistent with providing the largest incentives for such investments. However, in Section 4.1 we discuss how our results would extend if the goal were to maximize the payoff of the unviable types or the total payoff of all types.

**Comments** In our model, the evidence exchanged is instrumental to the project: the project cannot be undertaken unless all evidence is exchanged. Moreover, once information is exchanged, cooperation happens automatically. We find this assumption reasonable in a variety of settings. For example, the research project cannot be completed unless all ideas are put down in the paper, and as soon as they are written down, the decision on whether to submit the project for publication or shelve it is trivial. However, in a variety of settings the decision about

cooperation may be undertaken even before all information is exchanged. For example, consider two firms entertaining a merger. They exchange proprietary information that includes financial statements, contracts with other firms, and revenues; they visit each other's factories; and most generally, they "open the books" to each other. However, they can commit to the merger even before they exchange all information, that is, before they are certain that the opponent is viable and the merger will result in synergies. Similarly, the information that criminals exchange may only serve to signal that they are not undercover cops, and they may engage in a successful criminal venture even without exchanging this information. Situations like these can be easily mapped into our model with one modification: one has to endogenize the decision of cooperation and hence the amount of information exchanged. It is straightforward to show, however, that if the privacy concerns are smaller than the disutility from engaging in a venture with an unviable opponent, players will exchange all information before they decide to cooperate. In such cases all our results will continue to hold.

Next, a few comments on our assumptions are in place. First, if players had no privacy concerns, all projects could be undertaken with only a two-round evidence exchange—player 1 would reveal all her evidence in round 1, and then player 2 would reveal all her evidence in round 2. All equilibria leading to the project would deliver the same welfare.

Second, in the current model only the amount of evidence revealed, and not its order, matters. Such a modeling assumption is clearly appropriate if there is only one feasible order of evidence exchange (e.g., subsequent lemmas feed on the previous ones) or if a player cannot distinguish ex ante between different pieces of the opponent's evidence, and hence is unable to require them in any particular order. However, our model is more general: If each player can require the opponent to reveal her evidence in a particular order, then once this order is agreed upon, it results in some  $F$ ,  $h$ , and  $g$ , and our analysis follows. We will discuss this further in Section 3.2.1.

And finally, the problem of privacy could be easily solved if an expert benevolent mediator were available. However, in many circumstances a mediator with expertise sufficient to judge the viability of the presented evidence is unlikely to be benevolent.

## 3 Analysis

### 3.1 Preliminaries

Our first proposition states that we can divide all equilibria into two categories.

**Proposition 3.1** *In any equilibrium, either*

- a. the players cooperate on the project with probability 1 when  $\tau_1 = \tau_2 = V$ , or*
- b. viable types reveal no evidence.*

We delegate most proofs to the Appendix, but provide the proof of Proposition 3.1 here as it is short and conveys the intuition.

**Proof:** Since we consider pure strategy equilibria, if the project is undertaken with some probability when players are viable, this probability must be 1. Suppose then that there exists an equilibrium in which the project is never undertaken. By Assumption 2.1, there is a last round  $T(\sigma_j)$  at which the viable type of player  $i$  reveals new evidence. Suppose without loss of generality that  $T(\sigma_j) < T(\sigma_i)$ . Then at  $T(\sigma_i)$ , player  $i$  knows that if she reveals the piece of evidence prescribed by the equilibrium, she suffers a disutility from that, and in return she can at most receive evidence from the unviable type, which is not valuable. Hence, not revealing any evidence at  $T$  is a profitable deviation. ■

We will call the equilibria from part (a) *cooperating* and the equilibria from part (b) *non-cooperating*.

**Lemma 3.2** *Any cooperating equilibrium strictly Pareto dominates any non-cooperating equilibrium.*

The intuition for Lemma 3.2 is simple: since viable players always have the option to reveal no evidence, they must be better off in any equilibrium in which they willingly reveal evidence.

Given Lemma 3.2, we will henceforth focus on cooperating equilibria. The following proposition outlines the most important aspects of all such equilibria.

**Proposition 3.3** *In any cooperating pure strategy PBE there exists  $T > 0$  and a sequence  $\{\bar{N}_t\}_{t=1}^T$  with  $\bar{N}_{T-1} = \bar{N}_T = 1$  such that in each  $t$ ,*

- a. after any  $H_{t-1}$  with  $\{N_1^1, N_2^2, \dots, N_{t-1}^i\} = \{\bar{N}_1, \bar{N}_2, \dots, \bar{N}_{t-1}\}$  and  $s_{t-2}^j = s_{t-1}^i = \emptyset$ , all viable types of  $j$  and all unviable types of  $j$  with  $K^j > \bar{N}_t$  reveal  $\bar{N}_t$ ;*
- b. if  $N_{t-1}^i < \bar{N}_{t-1}$ , then the viable type  $j$  reveals  $N_t^j = \bar{N}_{t-2}$ ;*
- c. if  $s_{t-1}^i = U$ , then the viable type  $j$  reveals no new evidence in the game.*

*If  $pv - h(1) \geq 0$ , then a cooperating equilibrium exists and any sequence  $\{\bar{N}_t\}_{t=1}^T$  with  $\bar{N}_{T-1} = \bar{N}_T = 1$  can be supported as a cooperating equilibrium.*

Proposition 3.3 fully characterizes behavior on the equilibrium path: players adhere to the prescribed sequence of evidence revelation as long as no one has deviated (part *a*). As soon as one player deviates (unless she deviates to revealing everything and proves that she is viable), the exchange of valuable evidence stops (parts *b* and *c*).<sup>4</sup> Note that the proposition does not characterize the strategies of the players off the equilibrium path, as those may vary across equilibria.

---

<sup>4</sup>To be precise, part *b* does not say that the exchange of valuable evidence stops completely after any deviation from  $\{\bar{N}_t\}_{t=1}^T$ , but in the appendix we show that it stops after all deviations from this sequence that happen on the equilibrium path.

It should not come as a surprise that many sequences  $\{\bar{N}_t\}_{t=1}^T$  can be supported as cooperating equilibria. This is simply because for a given sequence, adherence to this sequence can be enforced by off-equilibrium beliefs that consider every deviation (except possibly to  $N_t^i = 1$ ) as coming from the unviable type. The only constraints that  $\mathbf{N}_T$  must satisfy is that every viable player prefers to adhere to the sequence instead of walking away after some round  $t$  with the evidence received from the opponent. The condition  $pv - h(1) \geq 0$  assures that walking away is never profitable for a viable type no matter how we split the evidence.

From now on, we will identify each cooperating equilibrium with its corresponding sequence  $\mathbf{N}_T \stackrel{\text{def}}{=} \{\bar{N}_t\}_{t=1}^T$  of evidence revelation. Since we are interested in whether more incremental—and hence taking place in more rounds—communication is beneficial, we will restrict our attention to sequences  $\mathbf{N}_T$  for which  $\bar{N}_t \neq \bar{N}_{t-2}$  for all  $t \leq T$ . Such a restriction is without loss of generality and allows us to treat  $T$  as a measure of how incremental different sequences of evidence exchange are.

We conclude this section by deriving the payoffs of the viable players as a function of  $\mathbf{N}_T$ . Consider an equilibrium  $\mathbf{N}_T$ , and observe that if the viable type of player  $i$  faces another viable type, the two will exchange all evidence and receive  $v + g(1) - h(1)$  each. If, however, a viable type of player  $i$  faces an unviable type, she will stop revealing evidence—and hence stop incurring disutility from privacy losses—as soon as  $\bar{N}_t^j > K^j$ . This is because in such a round  $t$ , either she will receive the signal that her opponent is not viable, or her opponent will deviate from the prescribed sequence  $\mathbf{N}_T$ . The probability that the unviable opponent has  $K^j \in (\bar{N}_{t-2}, \bar{N}_t)$  is equal to  $F(\bar{N}_t) - F(\bar{N}_{t-2})$ . Thus, the expected utility of a viable player 1 is

$$E[u_1^V(\mathbf{N}_T)] = p(v + g(1) - h(1)) - (1 - p) \sum_{k=0}^K (F(\bar{N}_{2k+2}) - F(\bar{N}_{2k})) h(\bar{N}_{2k+1}), \quad (3)$$

where  $K = \frac{T-2}{2}$  if  $T$  is even and  $K = \frac{T-3}{2}$  if  $T$  is odd. The expected utility of a viable player 2 is derived similarly.

Equation 3 reveals that differences in payoffs across cooperating equilibria come only from differences in the expected disutilities from evidence revealed to unviable types. It will be convenient to denote this disutility

$$\Psi_1(\mathbf{N}_T) \stackrel{\text{def}}{=} \sum_{k=0}^K (F(\bar{N}_{2k+2}) - F(\bar{N}_{2k}))h(\bar{N}_{2k+1}), \quad (4)$$

and call it the *privacy leakage* of player 1. Similarly, denote by  $\Psi_2(\mathbf{N}_T)$  the privacy leakage of player 2. Hence, the equilibria that maximize the joint payoff of the viable types  $E[u_1^V(\mathbf{N}_T)] + E[u_2^V(\mathbf{N}_T)]$  are those that minimize the joint privacy leakage

$$\Psi_1(\mathbf{N}_T) + \Psi_2(\mathbf{N}_T) = \sum_{t=1}^{T-1} (F(\bar{N}_{t+1}) - F(\bar{N}_{t-1}))h(\bar{N}_t). \quad (5)$$

Note that in the shortest cooperating equilibrium,  $\mathbf{N}_2 = \{1, 1\}$ , a viable player 1 reveals all her evidence in the first round. Hence, in the second round, the viable type of player 2 knows which type she is facing. If she faces the unviable type she reveals nothing. If she faces the viable type, she reveals 1 unit of evidence and cooperation on the project is successful. Hence,  $\Psi_1(\mathbf{N}_2) = h(1) = 1$  and  $\Psi_2(\mathbf{N}_2) = 0$ . We call this equilibrium *simple*, and any other cooperating equilibrium *incremental*.

### 3.2 Simple Versus Incremental Evidence Exchange

Before we state the first result, it is useful to understand that, in our model, the value of evidence is two-fold. First, evidence has an *intrinsic* value—the actual content that makes it relevant for the success of the project—such as the description of the project design, the relevant computer code, or a proof. The more evidence a viable player reveals, the higher its intrinsic value. The intrinsic value of evidence is measured by  $h$ , and it counts as a loss for the viable player

who reveals it. However, the evidence revealed by a player also carries information about the type of this player: the more evidence a player reveals, the more likely the opponent is to believe that she is the viable type. This *extrinsic* value of evidence is measured by  $F$ . Whenever a player is called upon to reveal  $\bar{N}_t$ , the viable type of this player suffers from the intrinsic value lost, but the viable type of the other player benefits from the extrinsic value gained.

It turns out that what matters for the optimality of equilibrium exchange is the precise relationship between the intrinsic and extrinsic values of evidence. To summarize this relationship, it is convenient to use the following change of variables:

$$M \stackrel{\text{def}}{=} h(N). \tag{6}$$

$M$  measures the units of intrinsic value contained in an amount  $N$  of evidence. The expression  $F(h^{-1}(M))$  then measures the extrinsic value associated with  $M$  units of intrinsic value.

We are now ready to state the conditions under which the optimal equilibria are simple.

**Proposition 3.4** *Suppose that  $pv > h(1)$ .*

- a. If  $M = F(h^{-1}(M))$  for all  $M \in (0, 1)$ , then all cooperating equilibria deliver the same joint payoff to the viable types.*
- b. If  $M > F(h^{-1}(M))$  for all  $M$ , then the unique cooperating equilibrium that maximizes the joint payoff to the viable types is the simple one.<sup>5</sup>*

For intuition, suppose that at a certain stage of communication player 1 has revealed less evidence than her opponent. Who should reveal new evidence next? In (a),  $F(h^{-1}(M))$  is assumed to be linear in  $M$ , and so revealing an additional unit of intrinsic value always delivers one additional unit of extrinsic value. Hence,

---

<sup>5</sup>In fact, any sequence in which one player reveals all her evidence first is an optimal equilibrium, but recall that we are assuming that there are no “silent” stages.

it does not matter which player reveals new evidence next – the total gain/loss from the next piece of evidence will be the same. This implies that any mode of evidence exchange is optimal.

To understand the intuition behind part (b), it is easier to focus on the case in which  $F(h^{-1}(\cdot))$  is strictly convex, which is a sufficient condition for  $M > F(h^{-1}(M))$ . Under strict convexity, any unit of intrinsic value revealed by a player delivers less extrinsic value than each additional unit. This implies that a unit of intrinsic value revealed by the player who lags in the exchange carries less extrinsic value than a unit revealed by the player who is ahead. This means that as soon as player 1 reveals some evidence, it is optimal to ask her to reveal the rest of her evidence before player 2 speaks. This, in turn, implies simple evidence exchange.

When  $M > F(h^{-1}(M))$  but  $F(h^{-1}(\cdot))$  is locally concave, it is no longer true that at any stage of the exchange it is optimal to ask the player who is ahead to reveal the rest of her evidence. However, since the revelation of all evidence carries the same intrinsic and extrinsic value (by the normalization  $F(h^{-1}(1)) = 1$ ), and the revelation of  $N < 1$  units of evidence delivers less extrinsic value than its intrinsic value, the result still holds.

The intuition above immediately suggests that when  $F(h^{-1}(\cdot))$  is concave, incremental evidence exchange should be optimal. In this case, a unit of intrinsic value revealed by the player who lags in the exchange carries more extrinsic value than a unit revealed by the player who is ahead. Hence, optimality requires that in equilibrium the player who lags in evidence exchange is the one who reveals an additional unit of evidence. This implies that evidence should be exchanged in turns, and the pieces exchanged should be as small as possible. Proposition 3.5 below formalizes this intuition.

**Proposition 3.5** *Suppose that  $pv > h(1)$ .*

- a. If  $M < F(h^{-1}(M))$  for some  $M$ , then there exists an equilibrium that delivers a higher total payoff to the viable types than the simple equilibrium.*

b. Suppose that  $F(h^{-1}(\cdot))$  is strictly concave, and consider a cooperating equilibrium  $\mathbf{N}_T$ . Then there exists another cooperating equilibrium with  $T + 1$  rounds of communication that delivers a strictly higher total payoff to the viable types.

When  $F(h^{-1}(\cdot))$  is not globally concave, as in part (b), it is not always beneficial to let the player who lags reveal the next piece of evidence. Hence, it is not true that more incremental exchange is always better. However, when  $M < F(h^{-1}(M))$  for some  $M$ , it is definitely beneficial for both players to reveal  $M$  units of intrinsic value before exchanging any more evidence.

We would like to stress that the above propositions do not let us compare the payoffs between any two arbitrary equilibria with  $T > 2$  and  $T' > T$ . To see why, suppose that  $F(h^{-1}(\cdot))$  is strictly concave. Suppose that in both equilibria players reveal evidence in similarly sized pieces in each round. Then by the arguments presented above, the longer equilibrium should deliver a higher welfare. Suppose now that in the longer equilibrium with  $T' > T$ , player 1 reveals a large amount of evidence in the first round. Under concave  $F(h^{-1}(\cdot))$  this is suboptimal, and hence the longer equilibrium may deliver lower welfare.

### 3.2.1 Comments on the Shape of the Utility Function

We have established that the nature of the optimal evidence exchange depends crucially on the shape of  $F(h^{-1}(\cdot))$ . The shape of  $F(h^{-1}(\cdot))$  is an empirical question, but we would like to develop intuition for when it is likely to be convex or concave.

Let us start by assuming that  $F(K) = K$ . In such a case,  $F(h^{-1}(\cdot))$  is concave when  $h$  is convex. And we should expect  $h$  to be convex if the proprietary evidence is of little value unless a large quantity of it is obtained.

Suppose now that  $h$  is linear. Then,  $F(h^{-1}(\cdot))$  is concave if  $F$  is concave. Recall that  $K^i$  is interpreted as the smallest amount of evidence of  $i$  that allows

$j$  to verify that  $i$  is unviable. We expect  $F$  to be concave in environments in which most of the invalid evidence can be spotted quickly (small  $K^i$ ).

In certain applications, the order in which evidence is revealed may be fixed by the nature of evidence. For example, revealing a proof may require revealing its steps in a predetermined order, as the steps may build upon themselves. In other application, however, the players may have more flexibility as to the order in which evidence is presented. Our analysis suggests that in the latter case, they should order evidence in such a way that  $F(h^{-1}(\cdot))$  is concave: the first pieces of evidence should be the ones that are most likely to be possessed only by the viable types, but that are also less costly if leaked to an unviable type.

### 3.3 Properties of Optimal Incremental Evidence Exchange

Proposition 3.5 implies that when  $F(h^{-1}(\cdot))$  is strictly concave, the optimal mode of equilibrium exchange does not exist. For any number of rounds  $T$ , players would always benefit by splitting evidence even finer. Since the model is only an abstraction of actual situations, one can conjecture that in reality there is a limit to the number of rounds in which players can engage. The next proposition describes how players should split information if they are bound by  $T$  rounds.

**Proposition 3.6** *Fix  $T$  and suppose  $F(h^{-1}(\cdot))$  is differentiable and strictly concave. If  $(F(h^{-1}(\cdot)))'$  is (weakly) concave/convex, then among the equilibria of length  $T$ , in the one that maximizes the joint welfare of the viable types, the sequence  $\{h(\bar{N}_t) - h(\bar{N}_{t-1})\}_{t=1}^{T+1}$  is (weakly) increasing/decreasing.*

Hence, if  $(F(h^{-1}(\cdot)))'$  is concave, then the amount of intrinsic information revealed in each round should increase as the communication progresses. In this case it is as if players “build trust” in the initial rounds, and once this trust is built, they exchange evidence more freely. In the other case, players become more “cautious” towards the end.

The intuition for Proposition 3.6 is as follows. Recall that if  $F(h^{-1}(\cdot))$  is strictly concave, then the extrinsic value of an additional unit of intrinsic value revealed by the player who lags is higher than the one revealed by the player who leads. Hence, an equilibrium is optimal if in each round, the difference in the amount of evidence revealed by the players is small. For a finite  $T$ , however, there is a limit to how small this difference can be, but one can make it larger in some rounds and smaller in others. It is crucial then that the rounds in which this difference is high coincide with the rounds in which the difference between the extrinsic value of an additional unit of intrinsic value revealed by the lagging player and the leading player is the smallest. And when  $(F(h^{-1}(\cdot)))'$  is concave, this difference is the smallest at higher levels of  $h(\bar{N}_t)$ , that is, in later rounds. Hence, in later rounds players can split information less finely.

From Proposition 3.5, we know under what conditions optimal equilibria are incremental. One may ask, however, what is the welfare gain that players can achieve by playing the incremental equilibria. To answer this question, one needs to assume a functional form for  $h$ . The next proposition gives an example of the welfare gain that can be achieved for a relatively simple class of functions.

**Proposition 3.7** *Suppose  $h(N) = F(N)^\alpha$  (or equivalently,  $F(h^{-1}(M)) = M^{\frac{1}{\alpha}}$ ) for some real number  $\alpha > 1$ . Then in the optimal equilibrium with  $T$  rounds, it holds that  $\Psi^i(\mathbf{N}_T) + \Psi^j(\mathbf{N}_T) \rightarrow 2/(\alpha + 1)$  as  $T \rightarrow \infty$ .*

Hence, for higher  $\alpha$ —which correspond to a more concave  $F(h^{-1}(\cdot))$ —the welfare gain from incremental communication is higher. In the limit as  $\alpha \rightarrow \infty$ , evidence that has very little intrinsic value carries a lot of extrinsic value; hence, in this case one can avoid privacy leakage almost completely.

## 4 Extensions

### 4.1 Other Optimality Criteria

In this section we discuss how our results extend with other optimality criteria, namely the joint payoff of all types and the joint payoff of the unviable types.

Note first that when both players are viable or both players are unviable, the details of evidence exchange do not matter. This is because the viable types end up cooperating no matter what exchange protocol they follow, and the unviable types do not gain any valuable evidence. Hence, the mode of communication matters only if one player is viable and the other is not.

With probability  $p$ , player 2 is unviable. In this case, player 1's privacy leakage is described by the equation for  $\Psi_1(\cdot)$  (equation 4). Player 2's privacy leakage (which is a gain in this case) is derived solely from the evidence revealed by player 1; hence, her privacy leakage will be like in equation (4) but with  $(-g)$  in place of  $h$ . So the total privacy leakage will be like in (4) but with  $h(N) - g(N)$  in place of  $h(N)$ . One can derive the privacy leakage in the case when player 1 is unviable similarly. Summing them, we obtain that the total expected privacy leakage is

$$\sum_{t=1}^{T-1} (F(\bar{N}_{t+1}) - F(\bar{N}_{t-1})) (h(\bar{N}_t) - g(\bar{N}_t)). \quad (7)$$

Hence, maximizing total welfare of all types is equivalent to minimizing (7). By the same argument, maximizing the joint payoff of only the unviable types requires minimizing

$$\sum_{t=1}^{T-1} (F(\bar{N}_{t+1}) - F(\bar{N}_{t-1})) (-g(\bar{N}_t)). \quad (8)$$

Hence, all our results still hold with the caveat that the conditions in the propositions need to be placed on the function  $F(w^{-1}(\cdot))$  instead of  $F(h^{-1}(\cdot))$ , where  $w(N) \stackrel{\text{def}}{=} h(N) - g(N)$  when we maximize the total payoff and  $w(N) \stackrel{\text{def}}{=} -g(N)$  when we maximize the total payoff of the unviable types.

Note that evidence exchange is socially beneficial if  $h(N) - g(N) < 0$ . Using  $w(N) \stackrel{\text{def}}{=} h(N) - g(N)$ , our model can be used to characterize when incremental exchange is socially beneficial in such settings.

## 4.2 All Types Have Privacy Concerns

So far we have assumed that the unviable types do not have privacy concerns. Even though this is likely to hold in most of our examples, one can easily imagine situations in which this is not true. The following provides one such example.

**Example 4.1** *Two firms are contemplating a merger based on the perceived potential synergies. To complete the merger, they have to share all private information; e.g., their financial statements, accounting procedures, initiated investments, corporate culture. Each firm  $i$  privately knows whether it satisfies conditions for the synergies to be realized, and if it does not, then this becomes apparent to firm  $j$  after  $i$  reveals  $K^i$  amount of information. In this example, both firms may be concerned about their privacy independently of their type: In the event of the merger not occurring, both firms can use the acquired information in the marketplace.*

In this section, we discuss how our results would change if we allowed all types to have privacy concerns.

Let  $h_{\tau_i}(N^i)$  be the disutility that player  $i$  suffers if she reveals  $N^i$ , and let  $g_{\tau_j}(N^j)$  be the utility that player  $i$  receives if she obtains  $N^j$  from her opponent. Hence, the payoffs may depend on whether the source of evidence is  $U$  or  $V$ :

$$u_i^V = v \mathbf{1}_{(\tau_j=V, N^1=N^2=1)} - h_V(N^i) + g_V(N^j) \mathbf{1}_{(\tau_j=V)} + g_U(N^j) \mathbf{1}_{(\tau_j=U)}, \quad (9)$$

and

$$u_i^U = -h_U(N^i) + g_V(N^j) \mathbf{1}_{(\tau_j=V)} + g_U(N^j) \mathbf{1}_{(\tau_j=U)}. \quad (10)$$

As before, all functions are continuous, take value 0 at 0,  $g_V$  and  $h_V$  are strictly increasing, and we normalize  $h_V(1) = 1$ . Until now, we were assuming that

$h_U \equiv g_U \equiv 0$ .

It should be clear that Proposition 3.1 and Lemma 3.2 still hold. We show in the appendix that the behavior outlined in parts *a*, *b*, and *c* of Proposition 3.3 holds with possibly one exception: the unviable types do not have to adhere to the sequence  $\mathbf{N}_T$  (the second part of *a*). They follow the sequence only up to a certain round, after which they reveal no new evidence. This is because the unviable types now have privacy concerns as well, and therefore adhering to the equilibrium sequence may result in too much privacy loss for them to be optimal.

In the presence of privacy concerns for the unviable types, making evidence exchange more incremental has an additional effect: If the difference between the amount of information revealed in each round is small enough, some of the unviable types may find it optimal to stop revealing evidence early in the game (possibly in their first round). This decreases the information leakage of the viable types; hence, it should be beneficial. Indeed, Proposition 4.2 below states that when the viable types have any privacy concerns, the simple equilibrium is *never* optimal.

**Proposition 4.2** *Suppose that  $h_U$  is strictly increasing, and the simple equilibrium exists. Then there exists  $\bar{N}_1 \in (0, 1)$  and  $\bar{N}_2 \in (0, 1)$  such that  $\mathbf{N}_4 = \{\bar{N}_1, \bar{N}_2, 1, 1\}$  strictly Pareto dominates (in terms of the payoffs of the viable types) the simple equilibrium.*

Note that if  $\bar{N}_1$  is sufficiently large and  $\bar{N}_2$  sufficiently small, the unviable types of player 1 and 2 do not reveal any evidence in the game, as the privacy loss from revealing  $\bar{N}_1$  and  $\bar{N}_2$ , respectively, is larger than the possible evidence gain from  $\bar{N}_2$  and  $1 - \bar{N}_1$ , respectively. Hence, the viable type of player 2 knows the type of her opponent already in the second round. As a result, she reveals no evidence to the unviable opponent, obtaining the same privacy leakage 0 as in the simple equilibrium. Similarly, the viable type of player 1 knows the type of her opponent in the third round; hence, her privacy leakage is only  $h_V(\bar{N}_1)$ . And

this is strictly less than the privacy leakage in the simple equilibrium, namely  $\Psi_1(\mathbf{N}_2) = h_V(1)$ .

Let us call all equilibria that discourage the unviable types from evidence exchange *screening*. The crucial feature of the screening equilibrium  $\mathbf{N}_4$  is that the difference between what a player has to reveal in a given round and what he expects to receive in the next round is small enough that the benefit does not outweigh the privacy loss. This suggests that by splitting information finer, one can discourage the unviable types from participating using a lower  $\bar{N}_1$ , and hence achieving a lower privacy leakage  $h_V(\bar{N}_1)$ . The following Proposition 4.3 confirms this intuition.

**Proposition 4.3** *Let  $h$  be a strictly increasing and continuous function with  $h(0) = 0$ , and let  $h_U(N) = bh(N)$  for some  $b > 0$ . Then if  $\mathbf{N}_T$  is a screening equilibrium, then there exists  $\mathbf{N}_{T+1}$  that is also a screening equilibrium with  $N_1(\mathbf{N}_{T+1}) < N_1(\mathbf{N}_T)$ . Moreover, for each fixed  $T$ , the lowest  $N_1$  in any screening equilibrium is a decreasing function of  $b$ .*

The analysis above sheds light on the shape of the optimal equilibria. Suppose first that  $g_U(N) = 0$  for all  $N$ . That is, the viable types do not benefit from the evidence obtained from the unviable types. In this case, screening is beneficial. When  $F(h_V^{-1}(\cdot))$  is convex and  $h_U(1)$  is small, the optimal equilibrium is a screening equilibrium with as many rounds as possible. When  $F(h_V^{-1}(\cdot))$  is concave, then the optimal equilibrium is incremental either for the reasons outlined in the previous sections or because it discourages the unviable types from evidence exchange. It is not straightforward, however, to compare the screening equilibria with the nonscreening ones in this case, as the non screening equilibria may admit a lower  $\bar{N}_1$ .

The same logic favors incremental equilibria when  $g_U(N)$  is strictly increasing but small. When  $g_U(1)$  becomes large, however, the analysis becomes more complicated. This is because the viable types may prefer the unviable types to

reveal as much information as possible. Hence, when both  $g_U$  and  $h_U$  become large, there are two competing effects: one can discourage the unviable types as the expense of a smaller  $\bar{N}_1$ , but discouraging them becomes less valuable. Which effect dominates is left for future research. In Example 4.4 below, however, we demonstrate that at the other extreme when  $g_U \equiv g_V \equiv h_U \equiv h_V$ , the screening effect dominates and incremental exchange is optimal.

**Example 4.4** *Suppose  $g_U \equiv g_V \equiv h_U \equiv h_V$ . There exists  $T_0$  such that the following holds for all  $T \geq T_0$ : There exists a cooperating equilibrium with a sequence  $\mathbf{N}_T$  in which the viable types reveal all their evidence according to  $\mathbf{N}_T$  as long as their opponent does the same, but unviable types do not reveal any evidence. Moreover, for these equilibria,  $\Psi_2(\mathbf{N}_T) = 0$ , and as  $T \rightarrow \infty$ ,  $\Psi_1(\mathbf{N}_T) \rightarrow 0$ . Hence, one can construct an equilibrium in which the payoffs of the viable types are arbitrarily close to the payoffs they would obtain if they suffered from no privacy concerns.*

### 4.3 Discounting

Our interpretation of the real-world communication is that exchanging evidence always takes the same amount of time, independent of whether one player reveals all evidence first or players split evidence into smaller pieces that they then exchange in an alternating fashion. For that reason, we do not incorporate discounting into our model. However, if alternating evidence exchange requires more time or carries some other cost (e.g., cost of attention), then players have an incentive to minimize the number of rounds. In this case, the characterization of the optimal equilibria is less straightforward, but the general insights hold. Since discounting incentivizes players to have as few rounds as possible, the simple equilibria are still optimal when  $F(h^{-1}(\cdot))$  is convex. When  $F(h^{-1}(\cdot))$  is sufficiently concave, the optimal equilibria are still incremental, but unlike before, for each  $F(h^{-1}(\cdot))$ , there exists a finite optimal number of rounds  $T$ .

Discounting also affects the analysis of Section 4.2. The equilibrium  $N_4$  from Proposition 4.2 is still an equilibrium, but any other incremental equilibrium that discourages all unviable types from exchange is not: Once both players are revealed viable, they have a strong incentive to reveal all remaining evidence at once. In such a case, one can still construct equilibria in which some unviable types drop out in the first two rounds, but one has to leave some unviable types with large  $K^i$  in the game to incentivize the viable types to reveal evidence incrementally.

## 5 Conclusions

Our results indicate that optimal communication can be simple or incremental, depending on the order in which evidence is revealed. However, when players can decide on the order or when all players have privacy concerns, incremental evidence exchange dominates. Our paper sheds light on possible drivers behind the fact that, in practice, we often observe incremental communication.

## Appendix

### A Proofs from Section 3.1

#### A.1 Proof of Lemma 3.2

If there is no cooperating equilibrium, then the lemma is vacuously true. Suppose then that there exists a cooperating equilibrium that is weakly Pareto dominated by some non-cooperating equilibrium. Since by Proposition 3.1, in the non-cooperating equilibria no valuable evidence is exchanged, and so the payoff of each viable type is 0. Hence, if a cooperating equilibrium is dominated by a non-cooperating one, then the payoffs in the former must be  $E[u_1^V] \leq 0$  and  $E[u_2^V] \leq 0$ . First, it cannot be that  $E[u_i^V] < 0$ , as then player  $i$  would prefer not to reveal

her evidence at all and obtain at least 0. Second, in any cooperating equilibrium  $E[u_2^V] > 0$ , as player 2 has an option to walk away from communication after receiving the first round of evidence, which with probability  $p$  is valuable. Hence, her expected payoff in the equilibrium must be strictly positive. Actually, the lemma is true even if we consider joint payoff of all types (as in Section 4.1), as the unviable types can only benefit from evidence exchange.

## A.2 Proof of Proposition 3.3

In what follows,  $p_t^i$  denotes the posterior belief held at the end of period  $t$  by  $j$  about  $i$  being the viable type. Let  $\sigma_1$  and  $\sigma_2$  denote the pure strategies of the viable type of players 1 and 2 in a cooperating equilibrium. Let  $\{\bar{N}_t^1\}_{t \in \{1,3,\dots\}}$  and  $\{\bar{N}_t^2\}_{t \in \{2,4,\dots\}}$  be the amount of information that the viable types of players 1 and 2 reveal on the equilibrium path if both players are viable and use  $\sigma_1$  and  $\sigma_2$ . By Assumption 2.1, viable types achieve cooperation in a finite number of rounds; hence, there exists  $t$  such that  $\bar{N}_{t-1}^i = \bar{N}_t^i = 1$ . Let  $T$  denote the smallest such  $t$ . By definition, in any cooperating equilibrium at any  $t$ , the viable type of the speaking player  $i$  must adhere to  $\{\bar{N}_t^1\}_{t \in \{1,3,\dots\}}$  and  $\{\bar{N}_t^2\}_{t \in \{2,4,\dots\}}$  as long as the opponent has adhered to it so far and  $s_{t-2}^i = s_{t-1}^j = \emptyset$ . Setting  $\{\bar{N}_t\}_{t=1}^T = \{\bar{N}_1^1, \bar{N}_2^2, \bar{N}_3^1, \dots\}$  proves the behavior of the viable types outlined in part (a).

Step A: If in some  $t$ ,  $p_t^i = 0$ , then the viable type of  $j$  reveals no evidence in  $t+1$ . This is straightforward as given such beliefs,  $j$  expects only privacy leakage from continuing the evidence revelation.

Step B: If  $s_t^i = U$ , then the viable type  $j$  reveals no more evidence in the game. This comes directly from Bayes' rule and Step A. This proves part (c).

Step C: If in equilibrium the strategy of the unviable type of player  $i$  prescribes revealing  $N_t^i \neq \bar{N}_t^i$ , then if in round  $t$  the opponent  $j$  sees  $N_t^i$ , then she reveals no new evidence after  $t$  unless  $i$  reveals all her evidence and turns out to be viable. This follows directly from Bayes' rule and Step A.

Step D: To prove the behavior of the unviable types outlined in part (a), note that such types are indifferent between any amount of evidence they reveal, and strictly prefer to receive more evidence from the opponent. By Step C, if a strategy of  $j$  that does not adhere to  $\{\bar{N}_t\}_{t=1}^T$  is an equilibrium strategy, it results in the opponent  $i$  withholding her evidence in the first round in which  $j$  deviates from  $\{\bar{N}_t\}_{t=1}^T$ . Since adhering to  $\{\bar{N}_t\}_{t=1}^T$  as long as  $K^j > \bar{N}_t^j$  makes the viable opponent reveal more evidence in  $t+1$ , a strategy that does not adhere to  $\{\bar{N}_t\}_{t=1}^T$  as long as  $K^j > \bar{N}_t^j$  cannot be an equilibrium.

Step E: To prove part (b), note that if  $N_{t-1}^i < \bar{N}_t$  is on the equilibrium path, then part (b) follows from Step C. Suppose then that  $N_{t-1}^i$  is off the equilibrium path. Take the first round  $t-1$  in which  $N_{t-1}^i < \bar{N}_t$  is observed. Clearly, the beliefs of  $j$  at  $t-1$  can be arbitrary. If they are  $p_{t-1}^i = 0$ , then part (b) follows from Step A. Suppose then that  $p_{t-1}^i > 0$ , and that the viable  $j$ 's strategy is to continue revealing new evidence. Consider an unviable type with  $K_t^i \in (\bar{N}_{t-2}, \bar{N}_t)$ . If this type adheres to her equilibrium strategy, then in any case she obtains no more evidence from the opponent (either because of Step C or because her strategy requires revealing  $\bar{N}_t$  which results in  $s^i = U$ ). If she reveals  $N_t^i$  instead, she obtains some new evidence from the viable type of  $j$ . Hence, revealing  $N_t^i$  is profitable, which contradicts the assumption that  $N_t^i$  was off the equilibrium path.

*Last claim*

We will show now that if  $pv \geq h(1)$ , then any sequence  $\{\bar{N}_t\}_{t=1}^T$  with  $\bar{N}_{T-1} = \bar{N}_T = 1$  can be supported as the following cooperating equilibrium. Players adhere to the behavior outlined in parts (a), (b), and (c). Moreover, (d) as soon as  $N_t^i \neq \bar{N}_t$ , then the opponent (of either type) reveals no new evidence unless at some  $\tau > t$ ,  $i$  reveals  $N_\tau^i = 1$  and  $s_\tau^i = \emptyset$ ; and (e) if at any  $t$ ,  $N_t^i = 1$  and  $s_i = \emptyset$ , then  $j$  (of either type) reveals  $N_{\tau+1}^j = 1$ . There may be equilibria in which the behavior described in (d) and (e) does not hold, but it is straightforward to see that they will be outcome equivalent to the equilibrium outlined here. And since

we are proving existence, it is enough to prove the existence of one equilibrium.

The proof of the behavior described in parts (b) and (c), and the behavior of the unviable types described in (a) did not rely on the details of the sequence; hence, it will hold for any sequence. The behavior described in (d) and (e) is clearly optimal for any sequence. We will now show that each player has an incentive to adhere to the behavior outlined in part (a) for any sequence.

Consider one such sequence. Suppose that the players followed this sequence up to (but excluding) round  $t$ , and that a viable type of  $j$  moves in  $t$ . By Bayes' rule,  $p_{t-1}^i = \frac{p}{p+(1-p)(1-F(\bar{N}_{t-1}))}$ . If  $j$  adheres to the sequence, with probability  $p_{t-1}^i$  she will cooperate and with probability  $(1 - p_{t-1}^i)$  the evidence exchange will stop at some time  $\tau > t$ , the time at which the opponent reveals herself unviable. Hence, her expected payoff from adhering to the sequence is

$$p_{t-1}^i (v + g(1) - h(1)) - (1 - p_{t-1}^i) \sum_{k=\frac{t-1}{2}}^K \frac{F(\bar{N}_{2k+2}) - F(\bar{N}_{2k})}{1 - F(\bar{N}_{t-1})} h(\bar{N}_{2k+1}), \quad (11)$$

where again  $K = \frac{T-2}{2}$  if  $T$  is even and  $K = \frac{T-3}{2}$  if  $T$  is odd.

If  $j$  deviates in  $t$  to revealing all her evidence, then her expected payoff is

$$p_{t-1}^i (v + g(1) - h(1)) - (1 - p_{t-1}^i) 1,$$

which is clearly lower than her payoff from not deviating. If she deviates to anything else, then by (d), she expects no more evidence exchange. Hence her best deviation is to reveal no new evidence at  $t$ . In such a deviation, she suffers disutility  $h(\bar{N}_{t-2})$ , but with probability  $p_{t-1}^i$  she faces a viable type in which case she benefits from the evidence gained so far  $g(\bar{N}_{t-1})$ . Comparing this to (11) and using the formula for  $p_{t-1}^i$ , we obtain that she does not deviate if and only if

$$v + g(1) - h(1) \geq \frac{1-p}{p} \sum_{k=\frac{t-1}{2}}^K (F(\bar{N}_{2k+2}) - F(\bar{N}_{2k})) h(\bar{N}_{2k+1}) + g(\bar{N}_{t-1}) - \frac{h(\bar{N}_{t-2})}{p_t}.$$

Note that the summation on the right-hand side of the IC constraint is smaller than the privacy leakage in the entire game, which in turn we have shown is

smaller than 1. Using this, we obtain

$$\begin{aligned}
& \frac{1-p}{p} \sum_{k=\frac{t-1}{2}}^K (F(\bar{N}_{2k+2}) - F(\bar{N}_{2k})) h(\bar{N}_{2k+1}) + g(\bar{N}_{t-1}) - \frac{h(\bar{N}_{t-2})}{p_t} \\
& < \frac{1-p}{p} + g(\bar{N}_{t-1}) - \frac{h(\bar{N}_{t-2})}{p_t} \\
& < \frac{1-p}{p} h(1) + g(1) \\
& < v + g(1) - h(1),
\end{aligned}$$

where the last inequality is true if  $v > \frac{1}{p}h(1)$ . Hence, if  $v > \frac{1}{p}h(1)$ , the incentive compatibility constraint is satisfied.

### A.3 Proof of Proposition 3.4

As established before, the equilibrium that maximizes payoff of the viable types minimizes their total privacy leakage. Using the change of the variables introduced in (6), the formula for privacy leakage (5) becomes

$$\Psi_1(\mathbf{M}_T) + \Psi_2(\mathbf{M}_T) = \sum_{t=0}^{T-1} (F(h^{-1}(\bar{M}_{t+1})) - F(h^{-1}(\bar{M}_{t-1}))) \bar{M}_t. \quad (12)$$

If for all  $M$ , it holds that  $M \geq F(h^{-1}(M))$ , then

$$\begin{aligned}
\Psi_1(\mathbf{M}_T) + \Psi_2(\mathbf{M}_T) & \geq \sum_{t=0}^{T-1} (F(h^{-1}(\bar{M}_{t+1})) - F(h^{-1}(\bar{M}_{t-1}))) F(h^{-1}(\bar{M}_t)) \\
& = F(h^{-1}(\bar{M}_{T-1})) F(h^{-1}(\bar{M}_T)) = F(h^{-1}(h(1))) F(h^{-1}(h(1))) = 1,
\end{aligned}$$

where the inequality is strict if  $M > F(h^{-1}(M))$  for all  $M$  and it is an equality if  $M = F(h^{-1}(M))$  for all  $M$ . Since  $\Psi_1(\mathbf{N}_2) + \Psi_2(\mathbf{N}_2) = 1$ , the proposition follows.

### A.4 Proof of Proposition 3.5

To prove part (a), take  $M$  for which  $M < F(h^{-1}(M))$ , and let  $\bar{N}_1 = h^{-1}(M)$ . Then clearly  $h(\bar{N}_1) < F(\bar{N}_1)$ . Consider an incremental equilibrium  $\mathbf{N}_3 = \{\bar{N}_1, 1, 1\}$ .

We have

$$\Psi_1(\mathbf{N}_3) + \Psi_2(\mathbf{N}_3) = h(\bar{N}_1) + 1 - F(\bar{N}_1) < 1, \quad (13)$$

which is less than in the simple equilibrium  $\mathbf{N}_2$ .

To prove part (b), take an equilibrium  $\mathbf{N}_T$ , and consider a sequence  $\hat{\mathbf{N}}_{T+1}$  for which  $\hat{N}_t = \bar{N}_t$  for all  $t \in \{1, \dots, T-2\}$  and  $\hat{N}_{T-1} = x$ ,  $\hat{N}_T = 1$ , and  $\hat{N}_{T+1} = 1$ , where  $x \in (\bar{N}_{T-3}, 1)$ . Then

$$\begin{aligned} \Psi_1(\mathbf{N}_T) + \Psi_2(\mathbf{N}_T) - (\Psi_1(\hat{\mathbf{N}}_{T+1}) + \Psi_2(\hat{\mathbf{N}}_{T+1})) = \\ (1 - F(\bar{N}_{T-2})) + (1 - F(\bar{N}_{T-3})) h(\bar{N}_{T-2}) \\ - (1 - F(\bar{N}_{T-2})) h(x) - (F(x) - F(\bar{N}_{T-3})) h(\bar{N}_{T-2}) - (1 - F(x)). \end{aligned}$$

Thus, the privacy leakage in  $\mathbf{N}_T$  is strictly higher than in  $\hat{\mathbf{N}}_{T+1}$  if and only if

$$(1 - F(x))(1 - h(\bar{N}_{T-2})) < (1 - F(\bar{N}_{T-2}))(1 - h(x)),$$

which holds if and only if

$$\frac{1 - F(x)}{1 - h(x)} < \frac{1 - F(\bar{N}_{T-2})}{1 - h(\bar{N}_{T-2})}.$$

Using the change of variables  $\bar{M}_{T-2} \stackrel{\text{def}}{=} h(\bar{N}_{T-2})$  and  $z \stackrel{\text{def}}{=} h(x)$ , the above inequality can be rewritten as

$$\frac{1 - F(h^{-1}(z))}{1 - z} < \frac{1 - F(h^{-1}(\bar{M}_{T-2}))}{1 - \bar{M}_{T-2}}.$$

But when  $F(h^{-1}(\cdot))$  is strictly concave, one can find  $x \in (\max\{\bar{N}_{T-3}, \bar{N}_{T-2}\}, 1)$  and the corresponding  $z \in (\max\{\bar{M}_{T-3}, \bar{M}_{T-2}\}, 1)$  such that the above is satisfied. By Proposition 3.3, if  $pv > h(1)$ , then  $\hat{\mathbf{N}}_{T+1}$  is also an equilibrium, which completes the proof.

## A.5 Proof of Proposition 3.6

**Proof:** Let  $\{\bar{N}_t\}_{t=1}^T$  be the equilibrium of length  $T$  that minimizes the total

information leakage. Then by definition, we have

$$\{\bar{N}_t\}_{t=1}^{T-2} = \arg \min_{\{\bar{N}_t\}_{t=1}^{T-2}} \sum_{t=1}^{T-1} (F(\bar{N}_{t+1}) - F(\bar{N}_{t-1})) h(\bar{N}_t),$$

where  $\bar{N}_0 = 0$  and  $\bar{N}_T = \bar{N}_{T-1} = 1$ . From Proposition 3.5 we know that the solution is interior. Using the change of variables (6) in the above expression and differentiating with respect to  $M_t$ , we obtain that for each integer  $t$  such that  $1 \leq t \leq T - 2$  it must be the case that

$$\begin{aligned} & \frac{d}{dM_t} \left( \sum_{t=1}^T (F(h^{-1}(M_{t+1})) - F(h^{-1}(M_{t-1}))) M_t \right) \\ &= F(h^{-1}(M_{t+1})) - F(h^{-1}(M_{t-1})) + \frac{dF(h^{-1}(M_t))}{dM_t} (M_{t-1} - M_{t+1}) \\ &= 0. \end{aligned}$$

Hence,

$$\begin{aligned} \frac{dF(h^{-1}(M_t))}{dM_t} &= \frac{F(h^{-1}(M_{t+1})) - F(h^{-1}(M_{t-1}))}{M_{t+1} - M_{t-1}} \\ &= \frac{\int_{M_{t-1}}^{M_{t+1}} \frac{dF(h^{-1}(x))}{dx} dx}{M_{t+1} - M_{t-1}} \leq \frac{dF(h^{-1}(\frac{M_{t+1} + M_{t-1}}{2}))}{dM_t}, \end{aligned} \tag{14}$$

where the last inequality follows if the derivative of  $F(h^{-1}(\cdot))$  is weakly concave (the inequality is strict if the derivative of  $F(h^{-1}(\cdot))$  is strictly concave). When we compare the far left-hand side with the far right-hand side of the above inequality and use the fact that  $F(h^{-1}(\cdot))$  is strictly increasing, we obtain that

$$M_t \leq \frac{M_{t+1} + M_{t-1}}{2} \tag{15}$$

(with strict inequality if the derivative of  $F(h^{-1}(\cdot))$  is strictly concave). Thus, since this holds for every  $t \in \{1, \dots, T\}$  we get that  $\{M_t - M_{t-1}\}_{t=1}^T$  is (weakly) increasing. Using the change of variables (6), equation (15) implies the first part of the proposition. The proof for the convex case is analogous.  $\blacksquare$

## A.6 Proof of Proposition 3.7

Let  $\mathbf{N}_T = \{\bar{N}_t\}_{t=1}^T$  be the equilibrium of length  $T$  that minimizes the total information leakage. Then by definition, we have

$$\mathbf{N}_T = \arg \min_{\{\bar{N}_t\}_{t=1}^{T-2}} \sum_{t=1}^{T-1} F(\bar{N}_{t+1}) - F(\bar{N}_{t-1}) h(\bar{N}_t),$$

where  $\bar{N}_0 = 0$  and  $\bar{N}_T = \bar{N}_{T-1} = 1$ . We have

$$\begin{aligned} \Psi^i(\mathbf{N}_T) + \Psi^j(\mathbf{N}_T) &= \sum_{t=1}^{T-1} (F(\bar{N}_{t+1}) - F(\bar{N}_{t-1})) h(\bar{N}_t) \\ &= \sum_{t=1}^{T-1} (L_{t+1} - L_{t-1}) h(F^{-1}(L_t)) \\ &= \sum_{t=1}^{T-1} (L_{t+1} - L_{t-1}) L_t^\alpha \end{aligned}$$

using the change of variables  $L_t \stackrel{\text{def}}{=} F(\bar{N}_t)$  and the assumption that  $h(N) = F(N)^\alpha$ . Now, similarly to the proof of Proposition 3.6, after differentiating with respect to  $L_t$  the first-order conditions yield

$$L_{t+1} - L_{t-1} = \frac{h(F^{-1}(L_{t+1})) - h(F^{-1}(L_{t-1}))}{h(F^{-1}(L_t))'} = \frac{L_{t+1}^\alpha - L_{t-1}^\alpha}{\alpha L_t^{\alpha-1}}.$$

Thus,

$$\Psi^i(\mathbf{N}_T) + \Psi^j(\mathbf{N}_T) = \sum_{t=1}^{T-1} (L_{t+1} - L_{t-1}) L_t^\alpha \quad (16)$$

$$= \frac{1}{\alpha} \sum_{t=1}^{T-2} (L_{t+1}^\alpha L_t - L_{t-1}^\alpha L_t) + (1 - L_{T-2}). \quad (17)$$

Summing (16) with  $\alpha$  times (17) and observing that  $L_T = L_{T-1} = 1$  yields

$$(\alpha + 1)(\Psi^i(\mathbf{N}_T) + \Psi^j(\mathbf{N}_T)) = L_{T-2}^\alpha L_{T-1} + L_{T-1}^\alpha L_{T-2} + (\alpha + 1)(1 - L_{T-2}),$$

and so

$$\Psi^i(\mathbf{N}_T) + \Psi^j(\mathbf{N}_T) = \frac{2 + (L_{T-2}^\alpha - 1) + \alpha(1 - L_{T-2})}{\alpha + 1}.$$

The claim then follows from the fact that  $L_{T-2} \rightarrow 1$  as  $T \rightarrow \infty$  (by the proof of Proposition 3.5).

## B Proof from Section 4

### B.1 Proof of the version of Proposition 3.3 outlined in section 4.2

The proof of Proposition 3.3 depends on the unviable types not having privacy concerns only in two steps: A and D. When the unviable types have privacy concerns, Step A still holds but needs a different proof. Step D needs to be altered as the behavior of the unviable types is different than in Proposition 3.3. Below we present the new versions of these steps.

Step A: When  $p_t^i = 0$ , both players believe that evidence exchange will not lead to cooperation. Hence, at this point a player has an incentive to reveal new evidence only if the opponent is expected to reveal new evidence in return. But by the same argument as in the proof of Proposition 3.1, the player who is supposed to reveal new evidence last, has an incentive to deviate to revealing no new evidence. Hence, no evidence exchange can occur.

Step D: By Step C, if the equilibrium strategy of  $j$  does not adhere to  $\{\bar{N}_t\}_{t=1}^T$  in round  $t$ , it results in the opponent  $i$  withholding her evidence in  $t+1$ . Hence, it is better for  $j$  to either reveal no new information in this round, as this decreases her privacy leakage, or to adhere to the sequence  $\{\bar{N}_t\}_{t=1}^T$ .

### B.2 Proof of 4.2

Let  $\mathbf{N}_4 = \{\bar{N}_1, \bar{N}_2, 1, 1\}$  be such that

$$h_U(\bar{N}_1) \geq pg_V(\bar{N}_2) \tag{18}$$

$$h_U(\bar{N}_2) \geq g_V(1) - g_V(\bar{N}_1). \tag{19}$$

By continuity of  $g_V$  and  $h_U$ , it is possible to find  $\bar{N}_1$  close to 1 and  $\bar{N}_2$  close to 0, for which (18) and (19) are satisfied. Below we will show that there exists an equilibrium with the above  $\mathbf{N}_4$  in which the viable types of both players adhere to  $\mathbf{N}_4$ , and the unviable types of both players reveal no evidence. Hence, in this equilibrium the information leakage of player 1 is  $h_V(\bar{N}_1)$  and of player 2 is 0. Since in the simple equilibrium the information leakage of player 1 is  $h_V(1)$ , and player 2 is 0,  $\mathbf{N}_4$  strictly Pareto dominates the simple equilibrium.

Suppose first that both types of player 2 follow the strategies outlined in the previous paragraph. The unviable type of player 1 knows that by revealing  $\bar{N}_1$ , she suffers  $h_U(\bar{N}_1)$ , and gains  $g_V(\bar{N}_2)$  only if she faces the viable opponent. Hence, her *IC* is

$$0 \geq pg_V(\bar{N}_2) - h_U(\bar{N}_1),$$

which is satisfied by (18). Suppose now that both types of player 1 follow the strategies outlined in the previous paragraph. Then at  $t = 2$ , the unviable type of 2 knows the type of her opponent. If the opponent revealed  $\bar{N}_1$ , then player 2 benefits  $g_V(\bar{N}_1)$  if she does not reveal anything. If she pretends to be viable and reveals  $\bar{N}_2$ , then she suffers a privacy loss, but she will receive all information from the opponent. Hence, her *IC* is

$$g_V(\bar{N}_1) \geq g_V(1) - h_U(\bar{N}_2),$$

which again is satisfied by (19).

It remains to show that the viable types of both players have an incentive to adhere to  $\mathbf{N}_4$ , but it should be clear (and is straightforward to show) that whenever players have an incentive to reveal evidence in a simple equilibrium, then they do in  $\mathbf{N}_4$ .

### B.3 Proof of Proposition 4.3

Suppose  $\mathbf{N}_T$  is a screening equilibrium. This requires that every unviable type prefers to reveal nothing in the first round at which she speaks instead of planning

to reveal evidence until some round  $t$ . At  $t = 1$ , the unviable type of player 1 knows that with probability  $(1 - p)$  she faces an unviable type, in which case no evidence will be revealed in  $t = 2$ . With probability  $p$  she faces the viable type, in which case she can stay in the conversation until any  $t$  such that  $K^i > \bar{N}_t$  and receive  $\bar{N}_{t+1}$  from the opponent. Hence, for all odd  $t$ , the following IC constraint must be satisfied:

$$p (g_V (\bar{N}_{t+1}) - h_U (\bar{N}_t)) - (1 - p) h_U (\bar{N}_1) \leq 0. \quad (20)$$

In round 2, no player will reveal evidence if no evidence is disclosed in round 1. For the unviable player 2 not to disclose any evidence in round 2 after she observes  $\bar{N}_1$  in round 1, it must be that she prefers to walk away with  $\bar{N}_1$  instead of planning to follow  $\mathbf{N}_T$  until some  $t$ . That is, for all even  $t$ , the following IC constraint must be satisfied:

$$g_V (\bar{N}_{t+1}) - h_U (\bar{N}_t) \leq g_V (\bar{N}_1). \quad (21)$$

Step 1: If  $\mathbf{N}_T$  satisfies the IC constraints and at least one inequality is strict, then there exists another  $\mathbf{N}'_T$  that satisfies the IC constraints with inequality and  $\bar{N}'_1 < \bar{N}_1$ .

Take the first  $t$  at which the IC constraint is satisfied with strict inequality. If  $t = 1$ , then by continuity and strict monotonicity of  $h_U$  and  $g_V$  one can decrease  $\bar{N}_1$  and keep (20) satisfied. Suppose then  $t = 2$ . Then one can decrease  $\bar{N}_2$  and keep (21) satisfied. Since  $\bar{N}_2$  only enters the first and second period *IC*'s, we need to make sure that (20) for  $t = 1$  is satisfied. But decreasing  $\bar{N}_2$ , relaxes (20) for  $t = 1$ , so one can decrease  $\bar{N}_1$ . Repeating the same argument for any  $t$  proves the step.

Step 2: Suppose  $\mathbf{N}_T$  is a screening equilibrium. By Step 1, we can assume that  $\mathbf{N}_T$  satisfies the IC constraints with strict equality. Consider decreasing  $\bar{N}_1$  by  $\varepsilon_1$ . In order to keep the IC for  $t = 1$  satisfied, one needs to decrease  $\bar{N}_2$  by some  $\varepsilon_2$ , which by continuity of all functions is a continuous function of  $\varepsilon_1$ . Therefore,

to keep the IC for  $t = 2$  satisfied, one needs to decrease  $\bar{N}_2$  by  $\varepsilon_2$ , which again by continuity is a continuous function of  $\varepsilon_1$ . Continuing this logic until  $T - 1$ , we obtain that we need to decrease  $\bar{N}_{T-1}$  by some  $\varepsilon_{T-1}$ . So far, by constructions, all *IC* constraints are satisfied until  $T - 1$ . We need to make sure that the *IC* constraint at  $T - 1$  is satisfied as well. Suppose that  $T - 1$  is odd. Then the *IC* is

$$p (g_V (1) - h_U (\bar{N}_{T-1} - \varepsilon_{T-1} (\varepsilon_1))) - (1 - p) h_U (\bar{N}_1 - \varepsilon_1) \leq 0.$$

Since in  $\mathbf{N}_T$ ,  $\bar{N}_{T-1} = 1$ , and by continuity of  $\varepsilon_{T-1}$  with respect to  $\varepsilon_1$ , one can find  $\varepsilon_1 > 0$  that will make this *IC* satisfied. And by construction,,  $\bar{N}_1 - \varepsilon_1 < \bar{N}_1$ .

Step 3: It remains to show that the lowest  $N_1$  is a decreasing function of  $b$ . Since by Step 1, the lowest  $N_1$  is achieved in the equilibrium in which all constraints are binding, it remains to show that increasing  $b$ , relaxes all the constraints. But this is immediate from (20) and (21).

## B.4 Proof of Example 4.4

Define  $h \stackrel{\text{def}}{=} g_U \equiv g_V \equiv h_U \equiv h_V$ . Take  $T$  even (an analogous proof can be constructed for  $T$  odd) and a sequence  $\mathbf{N}_T$  with  $\bar{N}_T = \bar{N}_{T-1} = 1$  such that (20) and (21) are satisfied with equality for all  $t$ . Then by adding the left-hand sides and the right-hand sides of (20) and (21) for the corresponding  $ts$  one obtains that for any odd  $K$ ,

$$h(1) - h(\bar{N}_{T-K}) = \frac{K + 2 - 2p}{p} \cdot h(\bar{N}_1).$$

Hence, for  $K = T - 1$ , we get

$$\begin{aligned} h(1) - h(\bar{N}_1) &= \frac{T + 1 - 2p}{p} \cdot h(\bar{N}_1) \\ \Leftrightarrow \frac{p}{T + 1 - p} &= h(\bar{N}_1). \end{aligned}$$

By setting  $\bar{N}_1 = h^{-1} \left( \frac{p}{T+1-p} \right)$ , we obtain a sequence for which all IC constraints for the unviable types are satisfied, and the unviable types reveal no evidence.

If the viable types adhere to  $\mathbf{N}_T$ , then the privacy leakage is  $\Psi_1(\mathbf{N}_T) = h(\bar{N}_1)$  and  $\Psi_2(\mathbf{N}_T) = 0$ , and  $\lim_{T \rightarrow \infty} h(\bar{N}_1) = \lim_{T \rightarrow \infty} \frac{p}{T+1-p} = 0$ . Hence, it remains to show that the IC constraints of the viable types are satisfied.

At  $t = 1$ , if the viable type of player 1 reveals  $\bar{N}_1$ , with probability  $1 - p$  she faces the unviable type, and no more evidence is exchanged, and with probability  $p$  she faces the viable type and cooperates on the project. Hence, her IC constraint at  $t = 1$  is

$$0 \leq pv - (1 - p)h(\bar{N}_1). \quad (22)$$

In all other odd rounds she knows whether her opponent is viable. If she faces the unviable type, she clearly has an incentive to reveal no evidence. If she faces the viable type, she must prefer to continue exchanging evidence and end up with  $v + g_V(1) - h_V(1) = v$ , instead of walking away at some  $t$  with the evidence that the opponent revealed in  $t - 1$ . Hence, her IC constraint in round  $t$  is

$$v \geq h(\bar{N}_{t-1}) - h(\bar{N}_{t-2}). \quad (23)$$

In each round in which player 2 speaks, she knows the type of her opponent, hence, her IC constraint in all even rounds is identical to (23). By (20) and (21),  $h(N_{t-1}) - h(N_{t-2}) \leq \max\left\{\frac{1-p}{p}h(\bar{N}_1), h(\bar{N}_1)\right\}$ , and since  $\lim_{T \rightarrow \infty} h(\bar{N}_1) = 0$ , there exists  $T_0$ , such that for all  $T \geq T_0$ , (23) is satisfied. Similarly, (22) is also satisfied.

## References

- ADMATI, A. and PERRY, M. (1991). Joint projects without commitment. *The Review of Economic Studies*, **58** 259–276.
- AUGENBLICK, N. and BODOH-CREED, A. (2013). Conversations about type: Privacy, grammars and taboos. *Mimeo*.
- BARDSLEY, P., CLAUSEN, A. and TEAGUE, V. (2008). Cryptographic commitment and simultaneous exchange.

- BLUM, M. (1983). How to exchange (secret) keys. *ACM Transactions on Computer Systems (TOCS)*, **1** 175–193.
- COMPTE, O. and JEHIEL, P. (2003). Voluntary contributions to a joint project with asymmetric agents. *Journal of Economic Theory*, **112** 334–342.
- COMPTE, O. and JEHIEL, P. (2004). Gradualism in bargaining and contribution games. *Review of economic studies*, **71** 975–1000.
- CRAWFORD, V. and SOBEL, J. (1982). Strategic information transmission. *Econometrica: Journal of the Econometric Society* 1431–1451.
- DAMGÅRD, I. B. (1995). Practical and provably secure release of a secret and exchange of signatures. *Journal of Cryptology*, **8** 201–222.
- DZIUDA, W. (2011). Strategic argumentation. *Journal of Economic Theory*, **146** 1362–1397.
- GANGLMAIR, B. and TARANTINO, E. (2013). Conversation with secrets.
- GOLDWASSER, S., MICALI, S. and RACKOFF, C. (1989). The knowledge complexity of interactive proof systems. *SIAM Journal on Computing*, **18** 186–208.
- HÖRNER, J. and SKRZYPACZ, A. (2011). Selling information. *Cowles Foundation Discussion Paper No. 1743R*.
- KARTIK, N. (2009). Strategic communication with lying costs. *Review of Economic Studies*, **76** 1359–1395.
- LI, H., ROSEN, S. and SUEN, W. (2001). Conflicts and common interests in committees. *American Economic Review* 1478–1497.
- LOCKWOOD, B. and THOMAS, J. (2002). Gradualism and irreversibility. *Review of Economic Studies*, **69** 339–356.

- MARX, L. and MATTHEWS, S. (2000). Dynamic voluntary contribution to a public project. *Review of Economic Studies*, **67** 327–358.
- MILGROM, P. (1981). Good news and bad news: Representation theorems and applications. *The Bell Journal of Economics* 380–391.
- MILGROM, P. and ROBERTS, J. (1986). Relying on the information of interested parties. *The RAND Journal of Economics* 18–32.
- PITCHFORD, R. and SNYDER, C. (2004). A solution to the hold-up problem involving gradual investment. *Journal of Economic Theory*, **114** 88–103.
- SHIN, H. (1994). The burden of proof in a game of persuasion. *Journal of Economic Theory*, **64** 253–264.
- STEIN, J. C. (2008). Conversations among competitors. *The American Economic Review*, **98** 2150–62.
- WATSON, J. (1999). Starting small and renegotiation. *Journal of Economic Theory*, **85** 52–90.
- WATSON, J. (2002). Starting small and commitment. *Games and Economic Behavior*, **38** 176–199.
- YAO, A. C. (1982). Protocols for secure computations. In *Proceedings of the 23rd Annual Symposium on Foundations of Computer Science (FOCS)*. IEEE Computer Society, 160–164.