

# Can Society Function Without Ethical Agents? An Informational Perspective\*

Bruno Strulovici  
Northwestern University

June 14, 2022

[Click here for the most recent version.](#)

## Abstract

In order to function, society relies on many facts that must be learned through intermediaries with special expertise or access to information. This paper considers whether society can learn about such facts when intermediaries are devoid of ethical motives and act sequentially. The answer depends on the severity of *information attrition* affecting the amount of discoverable evidence about each fact. Information attrition is nonexistent in fields based on reproducible scientific evidence but can affect the evidence in criminal and corruption investigations. Applications to institution enforcement, social cohesion, scientific progress, and historical revisionism are discussed.

**Keywords:** Disinformation, Ethical Necessity, Fake News, Information Attrition, Institution Design, Mediated Learning, Social Learning.

---

\*This project has benefited from numerous conversations and comments, particularly from Alessandro Lizzeri, Alex Frankel, Ben Golub, Boli Xu, Doron Ravid, Francisco Poggi, Larry Samuelson, Laura Doval, Ludvig Sinander, Meg Meyer, Piotr Dworzak, Richard Zeckhauser, Ron Siegel, Simone Galperti, and Xavier Duran and audiences at various seminars. Early versions of this project were developed while I was visiting Microsoft Research New England and Harvard whose hospitalities are gratefully acknowledged.

# 1 Introduction

Journalism, science, branches of government, the military, and the corporate sector offer numerous examples of individuals who risk their life or their career to expose important facts. Journalists have been silenced, scientists have been condemned and ostracized, public officials and members of the military have seen their careers derailed, and corporate whistleblowers have been retaliated against for trying to learn and truthfully reveal the truth about facts of public interest.

To what extent do societies and institutions rely on such ethical agents to acquire information? Can institutions be designed with appropriately calibrated incentives to elicit the truth when agents lack intrinsic motives to seek and reveal the truth? Or are intrinsic motives necessary to achieve this objective? This paper introduces a framework to analyze these questions and provides results suggesting that the answer depends on the nature of the facts of interest and of the evidence available about these facts. The analysis rests on two concepts: *mediated learning* and *information attrition*, which are introduced next.

## 1.1 Mediated Learning

Many facts must be learned through agents with specific expertise or access to information. For example, the net benefits of a vaccine, the results of an election, the historicity of an event, and the manmade nature of climate change cannot be directly verified by the average citizen, and there is no public “epiphany” at which the truth is exogenously revealed to all. In these and many others instances, citizens must rely on intermediaries to learn anything about the fact of interest, a situation that we will call *mediated learning*.

To succeed, mediated learning relies on the investigative efforts and truthfulness of agents who can make false, misleading, or uninformed statements and are subject to biases, pressures, ambitions, cognitive and time costs and other considerations that may distort their behavior. In philosophy and various social sciences, it is common to consider agents whose rules of behavior prescribe, at least to some extent, to act truthfully.<sup>1</sup> By contrast, economists

---

<sup>1</sup>One branch of epistemic philosophy concerns the vulnerability of testimony, i.e., the fact that a speaker can lie, and resolves this vulnerability through behavioral assumptions driven by norms of truthfulness and principles of ethical behavior, such as Grice’s “cooperative principle” (Grice (1975)). Going further, Williams (2002) argues that agents must value truthfulness intrinsically in order to exchange information successfully

usually conceptualize truthful behavior as the result of properly calibrated incentives given to agents who are devoid of intrinsic motives for such behavior.

This paper examines the feasibility of mediated learning through agents who lack intrinsic motives to seek and reveal the truth. These agents will be called *non-ethical* to indicate that they are devoid of ethical motives to be truthful.<sup>2</sup> According to the definition used in this paper, non-ethical agents may care about material rewards and punishments and informational costs and biases, but not directly about the facts that they are supposed to investigate. For example, a prosecutor who gains from securing the conviction of a defendant and desires this outcome regardless of the defendant’s actual guilt is non-ethical.<sup>3</sup>

## 1.2 Information Attrition

The cornerstone of the analysis is the concept of *information attrition*, which captures the idea that information about a fact may be in limited, fragile, or diminishing supply. As a first illustration, consider the problem of determining whether a given athlete has used a banned substance during a competition. Drug tests typically rely on two blood or urine samples, “A” and “B”. If the first sample, “A”, is positive (suggestive of doping), the agency storing the second sample, “B”, is then asked to test this sample. At this point, the second sample is the only direct source of evidence left about whether the athlete benefited from a banned substance during the event: the amount of evidence concerning this fact has shrunk, and each test reduces the amount of information available in the case.

---

(see Section 3). Computer scientists consider the related problem of communication through potentially adversarial intermediaries, known in the literature as the *Byzantine generals* problem (Lamport, Shostak, and Pease (1982)). This literature assumes the existence of “loyal” generals who obediently follow the communication protocol set out by the planner. Finally, it is common in sociology to assume the existence of pro-social norms that facilitate truthful behavior. See Granovetter (2017) for a recent comparison of the paradigms in economics and sociology.

<sup>2</sup>The term “non-ethical” emphasizes the lack of an intrinsic disposition to seek and reveal the truth, as opposed to extrinsic motives to do so. In general, ethical agents may fail to be truthful and vice versa. For example, a host lying to a murderer about the location of a potential victim may be ethical. Conversely, an individual who always tells the truth regardless of consequences need not be ethical. This paper abstracts from this distinction and specifically defines non-ethical agents as having truth-independent preferences.

<sup>3</sup>Truth-dependent preferences do not imply ethical behavior. For instance, a prosecutor who wishes to convict innocent defendants and acquit guilty ones is unethical. This observation has no bearing on the main analysis of this paper, which focuses on truth-independent preferences. Mediated learning with ethical agents is discussed in Sections 3.1, 3.2, and 5.

Information attrition may be caused by various factors. The first one, outlined above, is that the process of learning about a fact sometimes requires transforming the evidence in a way that prevents its use by subsequent investigators: the signals are *disposable*.<sup>4</sup> Second, information attrition may be due to the exogenous degradation of evidence: as time goes by, evidence may deteriorate and become uninformative, for example due to poor storage conditions. Third, information attrition may be caused purposefully by individuals who tamper with or destroy some evidence, e.g., to dissimulate fraud.

Information attrition may also stem from social considerations. Experiments that require a large amount of resources (financial or human, as in the case of large-scale experiments) may become increasingly difficult to organize and unlikely to occur. In some cases, such as the determination of election results, the fact of interest is time sensitive and the number of experiments (e.g., vote recounts) is constrained by a deadline.

### 1.3 Incentive Unraveling with Information Attrition

To appreciate how information attrition affects mediated learning, let us first consider a situation in which the numbers of agents and of pieces of information are fixed, finite, and commonly known. Specifically, let us reconsider the sample-testing example, proceeding backward from the agency in charge of testing the second and last sample.<sup>5</sup> Since this agency possesses the only sample left, it can lie at no risk of being contradicted. If, for instance, the agency stands to gain publicity from incriminating the athlete, it can do so with impunity. If instead the agency benefits from exculpating the athlete (perhaps due to pressure from a sport governing organization), it can also do this without risking contradiction. And if the agency is indifferent between incriminating the athlete and absolving him, it has no incentive to incur the cost of running the test to begin with. In all these cases, the second agency's report is untethered to the truth.

---

<sup>4</sup>Other examples of disposable signals include human subjects in experimental psychology: once a person has been exposed to a particular experiment, this person is irremediably affected and no longer a good subject for the same experiment. More generally, it is a fundamental principle in physics that measurements are subject to the *observer effect*, whereby the observation of a physical system affects the state of the system. This principle is of particular importance for applications in quantum physics.

<sup>5</sup>While sample testing in sports is primarily used here for expositional purposes, unreliable handling of samples and the limited amount of samples can have severe consequences for athletes. For a recent example in boxing, see, e.g., <https://www.asahi.com/ajw/articles/14338804> and <https://www.japantimes.co.jp/sports/2021/06/15/more-sports/boxing-2/kazuto-ioka-false-positive-jbc/>.

Consider now the agency in charge of the first blood sample. Whatever report this agency produces, the second agency’s report will not be based on the truth, as shown in the previous paragraph. The first agency’s expected payoff is therefore independent of the content of its sample. Like the second agency, it has no incentive to report the true content of its sample, or even to learn this content, and incentives unravel. In this example, mediated learning is infeasible if agencies are non-ethical. When agencies are non-ethical, mediated learning is infeasible regardless of agencies’ material incentives.<sup>6</sup>

## 1.4 Reproducible Evidence and Incentive Design

Some investigations are not subject to information attrition. Consider the question of determining whether a mathematical proof is correct. The proof remains available for anyone to read no matter how many times it has been checked in the past. The unraveling argument no longer applies and this paper shows (Theorem 1) that mediated learning by non-ethical agents is feasible provided that agents’ material incentives are designed appropriately. The incentives that deliver successful mediated learning have an intuitive structure: an agent is rewarded if his report is vindicated by subsequent findings and punished otherwise, and the magnitude of the reward or punishment depends on how surprising the agent’s report was relative to the public belief about the state of the world before the agent’s report.<sup>7</sup>

Information attrition does not arise, either, in scientific inquiries that rely on reproducible experiments. For instance, no matter how many times physicists measure the weight of an electron, the experiment can be replicated more times. Facts based on reproducible evidence can, under the appropriate incentive structure, be learned through non-ethical agents.

As a result, the theory can explain why oaths of truthfulness are used in some contexts, in which information attrition is a concern, and not others. For instance, it may explain why

---

<sup>6</sup>In practice, the same laboratory is often in charge of storing and testing *both* samples, which may lead to even more severe incentive problems, as in the case of nationally organized doping schemes. If there are three or more agencies, the results are the same when the agencies proceed sequentially, regardless of the rule (e.g., majority rule) used to convict the athlete or incentives given to the agencies. Section 5.2 discusses the case in which agencies move simultaneously as well as other learning structures.

<sup>7</sup>Some of these features are similar to the socially optimal policy analyzed by Smith, Sørensen, and Tian (2021) in the context of herding models. In contrast to some results in the present paper (Theorems 2, 3, and 4), in herding models all equilibria are at least somewhat informative: uninformative cascades may occur, but only after—and, precisely, because—some information has been publicly revealed.

oaths are used in courts of law but not in “hard” sciences (see Section 5.3).

## 1.5 General Framework: Uncertain and Endogenous Information Attrition

To explore the effects of information attrition on mediated learning, this paper considers a framework in which investigators act sequentially over an infinite horizon and communicate through public reports, and information attrition may be uncertain and endogenous.

**Sequentiality.** The sequential nature of investigations is a key feature of the analysis.<sup>8</sup> This feature may be motivated by various activities, such as journalism and academic research, in which investigation outcomes often take the form of sequential publications. Likewise, police investigations typically proceed sequentially, with only one individual in charge of the investigation at any given time.

Another motivation for modeling the investigatory process as sequential and potentially unbounded is to avoid the existence of a terminal investigatory stage whose findings are definitive and can never be checked for distortions. For example, an enforcement or military agency sending investigators to learn about an event may be subject to biases and try to cover up or misrepresent the findings of its agents, unless the agency may with positive probability be itself subject to a higher level of monitoring. Sequentiality has the unique potential to hold every investigator accountable.<sup>9</sup>

**General Framework.** The general model expands the scope of the theory in three directions: (i) *exogenous uncertainty* about the amount of discoverable evidence, (ii) *sequential accountability*, and (iii) *endogenous uncertainty* concerning how evidence was handled by past investigators. These features are modeled as follows: (i) The supply of evidence available in a case may be unknown a priori, even to investigators; (ii) The number of potential investigators may be large and a priori unknown, even to current investigators. Unlike the second testing agency in the earlier example, all investigators can a priori be held accountable for their actions and declarations thanks to the existence of subsequent investigators;

---

<sup>8</sup>Other mediated learning structures, including non-sequential ones, are discussed in Section 5.2.

<sup>9</sup>Some institutions feature an investigatory level of last resort. Judicial systems in which a highest court makes unappealable decisions are an example. One may then be concerned that such institutions strongly rely for their success on the ethical standards of the individuals placed at the terminal investigatory level. Such a concern is consistent with the main results of this paper.

(iii) investigators may not know how past investigators have affected the evidence available. For example, consider the evidence left on a crime scene. Investigators typically do not know a priori the amount and nature of evidence left on the scene. Moreover, if an investigator inherits the case from a previous investigator, he may not know how diligent the previous investigator has been and, hence, what fraction of the evidence has already been discovered, lost, or destroyed. An investigator thus faces exogenous and endogenous uncertainty about the amount of discoverable evidence.

This paper provides conditions on the probability distribution of the supply of evidence under which mediated learning is feasible (Theorem 1), and under which it is not (Theorems 2, 3, and 4). When the necessary conditions are violated, mediated learning by non-ethical agents is impossible in a strong sense: even when (i) there is an unbounded sequence of potential investigators and (ii) investigators' incentives may be arbitrarily designed with full commitment of the designer, there does not exist *any* equilibrium in which at least one investigator provides an informative report with positive probability.

**Analytical Challenges.** It is easy to design incentives for which mediated learning always fails or, given some incentive structure, to construct an equilibrium in which mediated learning fails. The main challenge is to show that, in the presence of information attrition, mediated learning fails *for all incentive structures and all equilibria*, in the strong sense with probability one, *agents reveal no information at all about the fact*.

This result is established in a dynamic game in which the state variables in any given period include past reports and the public belief about the set of evidence that remains to be discovered. This set is not a priori bounded and thus lies in an infinite-dimensional space. Another difficulty is that the public belief about the number of pieces of evidence that remain to be discovered need not be monotonic: While discovering of a piece of evidence may suggest that the remaining supply of evidence is now smaller by one piece, this discovery could also reveal the tip of an iceberg of evidence, pointing to many more pieces to discover. Moreover, a discovery that contradicts past reports can indicate that previous investigators have lied and that whatever evidence they purported to have discovered (and thus “removed” from the supply of evidence) was in fact fabricated, resulting in a more optimistic belief about the amount of evidence that actually remains to discover.

An agent's belief about the remaining supply of evidence affects his incentives *directly*, through the probability that he will discover new evidence if he looks for it, and *indirectly*, through the probability that subsequent agents discover evidence that could be compared to

the first agent’s report. Unpacking this indirect channel, an agent cares about subsequent agents’ beliefs about the supply of evidence, their beliefs about subsequent agents’ beliefs, and so on. Moreover, an agent can manipulate subsequent agents’ beliefs about the supply of evidence by lying in his report.

Information attrition does not preclude the possibility of a long sequence of agents who work, or the possibility that many pieces evidence remain to be discovered, but it imposes a negative correlation between these two events. To exploit this negative correlation and show that an informative equilibrium cannot exist, one must compute bounds on agents’ beliefs along any given sequence. This is achieved by showing that the probability that an agent discovers evidence is probabilistically linked to the *expected* impact that this agent’s report has on subsequent agents’ *relevant* beliefs. This link is formalized by Proposition 2.

## 1.6 Outline

Section 2 introduces the formal model and the main results of the paper. In the baseline model, all agents must incur a cost to acquire information. This cost may be arbitrarily small, but strictly positive. A high-level presentation of the proof and its challenges is provided in Section 2.4. The model is then extended to allow for the existence of *witnesses*, who receive information for free and are subject to idiosyncratic, private biases, and *analysts*, who receive evidence for free but incur a cost from processing it.

Sections 3 and 4 discusses consequences and illustrations of information attrition. Section 5.1 describes the relation between the concepts of hard evidence, information attrition, and mediated learning. Section 5.2 considers alternative investigation structures. Section 5.3 discusses how to foster truth-dependent preferences among investigators. All results are proved in the Appendix.



## 2 Formal Analysis

This section contains the formal model and the main results of the paper. In the baseline model, agents must incur a cost to look for signals about the fact of interest. Theorem 1 shows that if (i) the supply of signals is unrestricted in a sense formalized by the theorem, and (ii) agents' rewards and punishments can be designed and taken to be sufficiently large, then there are instances in which agents can be incentivized to look for signals and truthfully report their findings. By contrast, Theorem 2 shows that mediated learning fails if the supply of signals is limited in a specific probabilistic sense, even when the cost of information acquisition is arbitrarily small.

The impossibility result is then extended to two other agency problems, in which signal acquisition is sometimes costless, but those agents who receive costless signals are subject to either private *reporting biases* (e.g., a witness holds a private grudge against a defendant) or to information *processing costs* (e.g., a laboratory receives a sample that it costly to analyze). In both cases, mediated learning fails even if private biases and information process costs are arbitrarily small (Theorems 3 and 4).

In reality, the same individual may be subject to both biases and informational costs, and the informational costs may include both information acquisition and an information processing costs.<sup>10</sup> For simplicity, the analysis assumes that each agent is subject to only one agency problem.

### 2.1 Baseline Model

A fact of interest,  $\omega \in \Omega$ , must be inferred from a sequence  $S = (s^1, \dots, s^{\tilde{K}})$  of signals, each of which takes values in some finite signal space  $\Sigma$ . The sequence  $S$  and its length  $\tilde{K} \leq \infty$  are stochastic. The joint distribution of  $(\omega, S)$  is arbitrary.<sup>11</sup>

---

<sup>10</sup>For example, an individual investigating an industrial accident may have to *look for* forensic evidence and then *process* this evidence to extract its meaning.

<sup>11</sup>One could impose some restrictions on this joint distribution. For example, one may impose that  $S$  reveals  $\omega$  perfectly or with some minimum level of precision. Such additional structure is not required for the general analysis of this paper. Alternatively, one could consider an information structure in which  $S$  always contains infinitely many signals, but the number of signals whose correlation with  $\omega$  and with one another exceeds any given positive threshold is finite and decreases with the threshold. The model in this paper adopts—for tractability—a different structure of dependence, in which discoverable signals are

In each round  $i \geq 1$ , a new agent arrives and makes two decisions: First, the agent privately chooses between seeking a signal (“working”) at cost  $c > 0$  and doing nothing (“shirking”). Second, the agent publicly sends a message  $m_i$  from some finite message space  $M$ . The agent can use mixed strategies.

Let  $S_i$  denote the sequence of signals that remain to discover at the beginning of round  $i$  (in particular,  $S_1 = S$ ). If the agent in round  $i$  (hereafter, “agent  $i$ ” or simply “ $i$ ”) works and  $S_i \neq \emptyset$ , then  $i$  discovers some element of  $S_i$  with probability  $\lambda \in (0, 1]$ .<sup>12</sup> The discovered signal is denoted  $s_i$ .<sup>13</sup> No constraint is imposed on how likely each element of  $S_i$  is of being discovered by agent  $i$ . For instance,  $i$  could always be discovering the first element of  $S_i$ . Alternatively, this likelihood of each signal could depend arbitrarily on  $i$ ’s identity, on  $S_i$ , and on past messages.<sup>14</sup> With probability  $1 - \lambda$ ,  $i$  discovers no signal. If  $S_i$  is empty,  $i$  surely discovers no signal.

For simplicity, we do not model the possibility that agents destroy signals without seeing them, or that signals disintegrate exogenously. These additional forms of information attrition would only strengthen the paper’s impossibility results (Theorems 2, 3, and 4).<sup>15</sup> The extension is more complex when some agents, such as witnesses, can discover signals for free. This extension is analyzed explicitly in Sections 2.5 and 2.6.

After the information-seeking stage,  $i$  sends a report  $m_i$  whose distribution in  $\Delta(M)$  can depend arbitrarily on what (if anything)  $i$  has observed during the information-seeking stage and on the reports  $m_1^{i-1} = (m_1, \dots, m_{i-1})$  made by previous agents.

Entering round  $i + 1$ , we have  $S_{i+1} = S_i$  if  $i$  did not discover any signal. If  $i$  discovered a signal, then  $S_{i+1}$  is a subsequence of  $S_i$  with length  $|S_{i+1}| = |S_i| - 1$ .

---

arbitrarily correlated with  $\omega$  (and modeled as signals belonging to  $S$ ) but may be in finite supply.

<sup>12</sup>To simplify exposition, we assume that an agent discovers at most one signal, but what is called a “signal” in the paper could stand for multiple pieces of evidence. An earlier version of the paper allowed the possibility that each agent discovers multiple signals. It yielded the same results at the cost of more notation.

<sup>13</sup>The index notation distinguishes  $s_i$ , which is the signal discovered by agent  $i$  (when applicable) and  $s^j$  is the  $j^{\text{th}}$  signal in the sequence  $S$ .

<sup>14</sup>Moreover, Theorem 2 can be generalized to the case in which the probability  $\lambda$  of discovering a signal is nondecreasing in the length of  $S_i$ .

<sup>15</sup>The possibility result (Theorem 1), which allows a geometrically distributed number of signals, may be interpreted as allowing an exogenous decaying rate of evidence. Likewise, we do not model the possibility that signals be generated over time. If new signals can arrive and other can disappear, the feasibility of mediated learning will likely depend on the relative rates and dynamics of growth and disintegration.

Let  $m = (m_1, m_2, \dots)$  denote the sequence of reports made by all agents. The realized utility of agent  $i$  is given by

$$U_i = V_i(m, \omega) - c \mathbb{1}_{i \text{ works}} \quad (1)$$

where  $V_i$  takes values in some compact interval  $[-R, R]$ .

Agent  $i$  is *non-ethical* if the function  $V_i$  is independent of  $\omega$ . This definition captures the idea that  $i$  does not care directly about the truth.<sup>16</sup> If  $i$  is non-ethical  $V_i$  may be defined on the restricted domain  $M^{\mathbb{N}}$ .  $V_i$  can depend arbitrarily on the entire sequence  $m = (m_1, m_2, \dots)$ . In particular, the impossibility results presented in this paper hold regardless of whether  $V_i$  is an exogenously given utility function or one that is specifically designed (or, at least, influenced) by a regulator or social planner.

To illustrate the various forms that  $V_i$  may take, note that  $i$  could be punished if his report is contradicted by subsequent investigators or rewarded if his report differs from past investigators' (as in a journalistic scoop).  $V_i$  may aggregate a discounted stream of rewards and punishments.<sup>17</sup> The general formulation captures situations in which  $i$ 's utility is affected by reports indirectly through the actions that these reports trigger. For example,  $i$  could be a prosecutor in a trial, whose outcome  $a(m) \in \{\text{'guilty'}, \text{'not guilty'}\}$  depends on the statements made by all agents involved in the case. If  $i$  is non-ethical, his utility may be modeled by  $V_i(m) = R \times \mathbb{1}\{a(m) = \text{'guilty'}\}$ , as in Landes' (1971) model of prosecutorial behavior. The model also encompassed the possibility that agent  $i$ 's realized utility could depend stochastically on other agents' reports. For example, suppose that the number of investigators is stochastic. We can model this by interrupting mediated learning at some stopping time  $\tau$ , in which case  $i$ 's realized utility depends only on  $(m_1, \dots, m_\tau)$  and  $i$  cares about his expected utility conditional on the information at the beginning of round  $i$ . Other variations, such as including a private type that affects  $i$ 's utility, are also easily encompassed by the model.<sup>18</sup>

Agents have a common prior about the distribution of  $S$ . The equilibrium concept is

---

<sup>16</sup>This concept is somewhat similar to the notion of being purely "extrinsically motivated" in Bénabou and Tirole (2003), and one could say that  $i$  is intrinsically motivated *by the truth* if  $V_i$  depends on  $\omega$ .

<sup>17</sup>For example, if  $i$  receives utility  $v_{i,j}(m_1, \dots, m_j)$  in round  $j \geq i$  and discounts future utility with some factor  $\delta < 1$ , then

$$V_i(m) = \sum_{j \geq i} \delta^{j-i} v_{i,j}(m_1, \dots, m_j).$$

<sup>18</sup>Such a private type is explicitly considered in Section 2.5.

(weak) Perfect Bayesian Equilibrium.<sup>19</sup>

## 2.2 Positive Result with Reproducible Evidence or Limited Attribution

For each  $k \geq 1$ , let  $F^k = \Pr(|S| \geq k)$  denote the prior probability that there are at least  $k$  signals to discover at the beginning of the investigation process.

**Definition 1** *An equilibrium is **informative** if at least one agent works with positive probability.*

When the survival function  $F^k$  decreases at most at a geometric rate, there are instances of the model and utility functions  $\{V_i\}_{i \in \mathbb{N}}$  for which mediated learning is feasible, as indicated by the following theorem, whose proof is in Appendix B.2.

**Theorem 1** *For any  $\rho \in (0, 1]$  and  $\lambda \in (0, 1]$ , there exist a joint distribution of  $(\omega, S)$  and utility functions  $\{V_i\}_{i \geq 1}$  such that the distribution of  $S$  satisfies  $F^k = \rho^{k-1}F^1$  for all  $k \geq 1$  and an informative equilibrium exists.*

Theorem 1 shows that even if agents are non-ethical and reports are public, mediate learning may be feasible provided that the supply of evidence is sufficiently large. Even if agents are non-ethical and can see past reports, it is possible to incentivize them to work as long as one can use sufficiently high rewards and/or punishments and these incentives are appropriately designed.

To illustrate Theorem 1, let us consider for now the case of in which  $\rho = 1$  and  $\lambda = 1$ . The assumption that  $\rho = 1$  captures reproducible evidence: it means that the supply of signals is unlimited.<sup>20</sup> The assumption that  $\lambda = 1$  means that every agent who works surely discovers a signal.

We construct a joint distribution over state and signals for which mediated learning can be made arbitrarily precise as long as the rewards and punishments are high enough.

---

<sup>19</sup>Impossibility results clearly hold for stronger equilibrium concepts, such as sequential equilibrium, and hold even for Bayes-Nash equilibria since off-path beliefs play no role in the analysis. Reciprocally, the equilibrium constructed to prove the positive result, Theorem 1, is sequential.

<sup>20</sup>The argument is almost identical if  $\rho < 1$ , as explained in the proof.

Thus suppose that  $S$  consists of infinitely many signals taking binary value, “ $H$ ” or “ $L$ ”. The signals are conditionally i.i.d.: there is an unknown state of the world  $\omega \in \{H, L\}$  such that each signal  $\tilde{s}$  satisfies  $\Pr(\tilde{s} = “H” | \omega = H) = \Pr(\tilde{s} = “L” | \omega = L) = \pi \in (1/2, 1)$ , and the signals are independently distributed conditional on  $\omega$ . Agents’ message space is chosen to be binary:  $m_i \in \{“H”, “L”\}$  for all  $i$ .

In the equilibrium that we consider, it is always in an agent’s interest to follow his signal if he acquired one. The relevant pure strategies are therefore: (i) to work at cost  $c > 0$  and report one’s signal ( $m_i = s_i$ ), (ii) to shirk and send message “ $H$ ” at no cost, and (iii) to shirk and send message “ $L$ ” at no cost.

Let  $p_1 = P(\omega = H)$  denote the prior about  $\omega$ . For any given equilibrium, let  $p_i$  denote the probability at the beginning of round  $i$  that  $\omega = H$  conditional on past reports  $m_1^{i-1}$  and  $\gamma_i$  denote the probability that  $i$  works conditional on past reports  $m_1^{i-1}$ .

**Proposition 1** *For any thresholds  $p_-$  and  $p_+$  such that  $0 < p_- < p_1 < p_+ < 1$ , there exist  $R > 0$ , utility functions  $\{V_i\}_{i \geq 1}$  taking values in  $[-R, R]$ , and thresholds  $\underline{p}, \bar{p}$  such that  $0 < \underline{p} < p_-$  and  $1 > \bar{p} > p_+$  for which the following strategy profile constitutes an equilibrium:  $\gamma_i = 1$  if  $p_i \in (\underline{p}, \bar{p})$  and  $\gamma_i = 0$  otherwise.*

The state  $\omega$  can thus be learned with arbitrary precision as long as the rewards and punishments used to incentivize agents can be taken to be high enough. In equilibrium, agents work with probability 1 until the posterior belief becomes extreme enough, at which point learning stops.

Since  $p_i$  is a martingale and each signal has the same level of informativeness,  $p_i$  must exit  $[\underline{p}, \bar{p}]$  with probability 1 in the candidate equilibrium. Incentives are provided as follows: if  $i$  reported “ $H$ ”, he gets a reward if  $\bar{p}$  is reached and a punishment if  $\underline{p}$  is reached, and vice versa. These rewards and punishment depends on the belief  $p_i$  before  $i$ ’s report. If  $p_i$  was very close to one of the boundaries and  $i$ ’s report takes the posterior away from this boundary,  $i$  gets a high reward if the belief process ends up exiting through the other boundary (a low probability event) and a very mild punishment if the belief process ends up crossing the nearby boundary.

It is beyond the scope of the present paper to describe the optimal way to incentivize learning. One would have to define an objective function for the designer and the optimal mechanism would depend on details of the distributions. Instead Theorem 1 shows that

mediated learning is compatible with a large class of instances of the model. The next sections show that what destroys mediated learning is the presence of information attrition.

### 2.3 First Impossibility Result

Say that an equilibrium is **uninformative** if it is not informative in the sense of Definition 1. According to that definition, uncovering any modicum of information about  $S$  (and, hence,  $\omega$ ) with positive probability, no matter how small, suffices to qualify an equilibrium as “informative.” The following theorem shows, however, that informative equilibria fail to exist when information is subject to attrition in a specific sense.

**Theorem 2** *For any parameters  $(R, c, \lambda)$ , there exist strictly positive thresholds  $\{\underline{F}^k\}_{k \geq 1}$  with the following property:*

*All equilibria are uninformative unless  $F^k \geq \underline{F}^k$  for all  $k \geq 1$ .*

**Corollary 1 (Bounded Support)** *For an informative equilibrium to exist, the support of  $|S|$  must therefore be unbounded.*

This corollary is an immediate consequence of Theorem 2: if the support of  $|S|$  is bounded by some constant  $K$ , then  $F^{K+1} = 0$ , which violates the threshold condition of Theorem 2 for all  $k \geq K + 1$  regardless of the parameters  $(R, c, \lambda)$  and utility functions  $\{V_i\}_{i \geq 1}$ .

Even if  $|S|$  has unbounded support, Theorem 2 implies that mediated learning is feasible only if the survival function  $F^k = \Pr(|S| \geq k)$  does not decrease too fast in  $k$ , i.e., so fast that  $F^k$  drops below the corresponding threshold  $\underline{F}^k$  for some  $k$ . When this violation occurs, even if a social planner could design agents’ payoff functions  $\{V_i\}_{i \geq 1}$  arbitrarily subject to the reward-punishment bound  $R$ , learning is guaranteed to fail.<sup>21</sup>

**Intuition:** To understand mediated learning failures, suppose for simplicity that there is a fixed number of discoverable signals. Successful learning requires that at least one agent discovers one of the signals, which can be incentivized only if a subsequent agent discovers (with positive probability) a second signal, which can be incentivized only if a

---

<sup>21</sup>It would be interesting to characterize the rate at which the cutoffs  $\{\underline{F}^k\}_{k \geq 1}$  go to zero. This question does not seem to have a simple answer. In particular, while Theorem 1 shows that there are settings in which geometric distribution of  $|S|$  are compatible with mediated learning, there is no indication that this condition is necessary.

third agent discovers a third signal with positive probability, and so on. Therefore, agents must reach with positive probability a round in which (i) the probability that a signal remains is arbitrarily low and (ii) the probability that some agent makes an informative report is positive. These conditions are incompatible: no agent wants to work when the expected benefit is arbitrarily close to zero. This shows that mediated learning must fail down this path, which causes incentives to unravel all the way back to the first agent and leads to a complete and global failure of mediated learning.<sup>22</sup> This intuition ignores important challenges, which are presented next together with a gist the rigorous proof. The formal proof is in Appendix A

## 2.4 Overview of the Proof and Challenges

For any  $i, k \geq 1$ , let  $F_i^k = \Pr(|S_i| \geq k \mid m_1^{i-1})$  denote the probability that there remain at least  $k$  signals to discover at the beginning of round  $i$  given past reports  $m_1^{i-1}$ . The prior probability  $F^k$  that the initial sequence  $S$  contains at least  $k$  signals satisfies  $F^k = F_1^k$ .

To understand how the proof works, suppose first that  $S$  contains at most one signal, i.e., that  $F_i^2 = 0$  for all  $i$ . We will show by contradiction that no agent ever works in equilibrium. When  $F_i^2 \equiv 0$ , the probability  $F_i^1$  that there is a signal to discover in round  $i$  is decreasing in  $i$  path by path.<sup>23</sup>

Agent  $i$  works only if two conditions hold: (i) the probability  $F_i^1$  that a signal remains is high enough—above a cutoff  $\underline{F}^1$  provided by Lemma 1, below—(ii) the probability that some agent  $j > i$  works after  $i$  has worked is high enough.

This creates a tension: On the one hand, the more likely agent  $i$  is to work, the larger the expected drop from  $F_i^1$  to  $F_{i+1}^1$ . On the other hand, an agent works only if subsequent agents also work with sufficient probability, which requires that  $F_j^1$  stay above  $\underline{F}^1$ . This dynamic

---

<sup>22</sup>This intuition differs from herding models (Bikhchandani, Hirshleifer, and Welch (1992) and Banerjee (1992)), in which learning failures (“cascades”) occur only after so much public information has been revealed that agents prefer to forgo their own signals. In the present model, the learning failure occurs from the first very agent. One way to visualize this feature is that cascades occurring in the future “ripple back” to the very beginning of the investigation sequence because agents’ payoffs in early rounds depend on the messages produced during the cascades.

<sup>23</sup>Intuitively, either  $i$  shirked, in which case his message is uninformative, or he worked, in which case he either found the only signal, and there is nothing left, or he found nothing, which makes agents more pessimistic that there is a signal left to be found.

is impossible as  $F_j^1$  must go down by a non trivial amount but becomes squeezed above  $\underline{F}^1$ .

To formalize this tension, consider any history up to round  $i$  and let  $M_i^+$  denote the set of messages  $m_i$  such that some  $j > i$  works with positive probability. We derive a formula that combines the following observations: (i) The probability that  $M_i^+$  occurs conditional on  $i$  working cannot be too small, because  $M_i^+$  contains all continuations for which  $i$ 's real findings matter and, hence, drives  $i$ 's incentive to work. (ii) Conditionally on  $M_i^+$ , the average drop in  $F_i^1$  is proportional to the probability that  $i$  works, because  $i$  affects beliefs only if he works: the likely  $i$  is to work ex ante, the larger the expected impact of his message on subsequent beliefs. These observations are captured by the following *Discovery-Belief (DB)* formula, derived in Proposition 2. Let  $\beta_i = \Pr_i(i \text{ discovers a signal})$  and suppose that  $F_i^2 = 0$ . Proposition 2 shows that there exists  $Q > 0$  s.t.

$$\beta_i \leq Q \frac{\mathbb{E}_i \left[ (F_i^1 - F_{i+1}^1(m_i)) \mathbb{1}_{m_i \in M_i^+} \right]}{F_i^1}.$$

The numerator captures  $i$ 's expected impact on beliefs *conditional on sending a message that yields an informative continuation equilibrium* ( $m_i \in M_i^+$ ). The denominator is bounded below by  $\underline{F}^1$ .

Now suppose by way of contradiction that there exists an informative equilibrium, and let  $\underline{F} = \inf\{F_j^1 : j \text{ works with positive proba}\}$  denote the infimum belief over all informative continuation equilibria. To build a contradiction, we choose an informative continuation equilibrium starting from some round  $i$  such that  $i$  works and  $F_i^1 \leq \underline{F} + \varepsilon$ . Such a continuation equilibrium must exist by definition of  $\underline{F}$ . Summing the DB formula over  $j > i$  yields

$$\Pr_i(\exists j > i \text{ who discovers a signal}) \leq \frac{Q}{\underline{F}^1} \mathbb{E}_i [(F_{i+1}^1 - F_J^1)] \quad (2)$$

where  $J$  (random, possibly infinite) indexes the last person who works.

We have  $F_J^1 \geq \underline{F}$  and, by monotonicity,

$$F_{i+1}^1 \leq F_i^1 \leq \underline{F} + \varepsilon. \quad (3)$$

From (2), the probability that some  $j > i$  discovers a signal must thus be of order  $\varepsilon$ , which is too small to incentivize  $i$  and yields the desired contradiction.  $\square$

## Challenges



The previous argument relies on two assumptions:  $F_i^2 = 0$  and  $F_i^1$  is decreasing. In general, beliefs  $F_i^k$  can be positive for all  $k$  and nonmonotonic in  $i$ . For example, agent  $i$  could send a message suggesting that previous agents have shirked and, hence, there are more signals to discover. Or  $i$  could send a message that is positively correlated with the existence of many other signal to discover, i.e., the equivalent of uncovering a “gold mine” of signals. The proof must address several challenges:

1. Suppose we wish to show that  $i$  never works if  $F_i^2$  is arbitrarily small (rather than exactly zero, as in the previous argument). Proposition 2 shows that the DB formula holds as long as  $F_i^2$  is sufficiently small relative to  $F_i^1$ . To apply Proposition 2 to all  $j > i$ , however, we need to show that  $F_j^2$  remains small for all  $j > i$ . Doob’s martingale inequality (Lemma 4) allows us to show that this event, denoted  $\mathcal{A}$  in the proof, is highly likely.

2. However, Doob’s martingale inequality holds with respect to a specific filtration, which is in the present case the filtration generated by past public reports at the beginning of each round. By contrast, agent  $i$ ’s incentives are driven by the probability of event  $\mathcal{A}$  *conditional on  $i$  working*. This probability can be very different from the unconditional probability at the beginning of round  $i$ , especially if the probability that  $i$  works is very small. This issue is addressed by Lemma 5.

3. The argument above selected  $i$  such that  $F_i^1$  was close to the infimum  $\underline{F}$ . When there are more than two signals, we would ideally like to choose  $i$  such that (i)  $F_i^2$  is arbitrarily small *and* (ii)  $F_i^1$  is close to  $\underline{F}$ , but we cannot impose both conditions. To address this, the proof uses a *boundary function* mapping  $f \mapsto \mathcal{F}^1(f)$ , which defines the infimum of beliefs  $F_i^1$  over informative continuation equilibria *such that  $F_i^2 \leq f$ , for  $f \in [0, 1]$* .

4. To exploit Equation (2), the argument above used that  $F_{i+1}^1$  was close to  $\underline{F}$  (see Equation (3)). This is generally false if  $F_i^1$  is nonmonotonic. This issue is addressed by Lemma 6.

## 2.5 Witnesses

We now introduce witnesses, who differ from the previous agents along two dimensions:

- They discover a signal for free.
- They are subject to (possibly, arbitrarily small) private preference shocks that affect which report they prefer to send.

In each round  $i \geq 1$ , agent  $i$  can be an investigator, identical to the agents of the baseline model, or a witness. Whether  $i$  is an investigator or a witness is public information.

If  $i$  is an investigator, the structure of  $i$ 's round, information, and utility is as in the baseline model. If  $i$  is a witness, he receives at no cost a signal  $s_i \in S_i$  at the beginning of the round, where  $S_i$  is the set of signals remaining at the end of round  $i - 1$ . In particular,  $i$  can be a witness only if  $S_i$  is nonempty. If  $i$  is a witness, the sequence  $S_{i+1}$  of available signals at the end of round  $i$  satisfies  $|S_{i+1}| = |S_i| - 1$ .

At the beginning of round  $i$ , the probability  $\varphi_i$  that  $i$  drawn to be a witness is equal to zero if  $S_i = \emptyset$ , and can take any value in  $[0, 1]$  otherwise.  $\varphi$  can depend on calendar time and on past reports  $m_1^{i-1}$ .<sup>24</sup>

In principle, the content of a witness' signal could be informative about the number of signals that remain to discover and, in particular, about whether information attrition is an issue for subsequent agents.<sup>25</sup> This would make it more difficult to state clear results linking the feasibility of mediated learning to the prior distribution of  $S$ .

We rule out this possibility and focus on the case in which a witness' signal never increases expectations about the total number of signals. This is achieved as follows: the sequence  $S$  of signals is obtained by, first, generating an infinite sequence  $S^\infty$  of signals, which may exhibit any arbitrary correlation between one another and, second, by truncating this sequence at some integer-valued random variable  $\tilde{K}$  that is independently distributed from  $S^\infty$ . Agents observe signals in the order of the sequence. Thus, writing  $S = (s^1, \dots, s^{\tilde{K}})$  (using superscripts to avoid confusion with the signals discovered in round  $i$ , which are denoted with subscripts), suppose that  $q_i$  signals have been uncovered by the beginning of round  $i$ , so that  $S_i = (s^{q_i+1}, s^{q_i+2}, \dots, s^{\tilde{K}})$ . If  $i$  discovers a signal  $s_i$ , then necessarily  $s_i = s^{q_i+1}$  and  $S_{i+1} = (s^{q_i+2}, \dots, s^{\tilde{K}})$ .

Moreover, we assume that  $\tilde{K}$  has an increasing hazard rate, i.e.,  $\Pr(\tilde{K} = k) / \Pr(\tilde{K} \geq k)$  is increasing in  $k$ .

**Assumption 1** (i) *The total number of signals  $\tilde{K}$  has an increasing hazard rate.* (ii) *For*

<sup>24</sup>For example, there could be a fixed (finite or infinite) subset  $\mathcal{N} \subset \mathbb{N}$  of rounds such that  $\varphi_i = \varphi^* \mathbb{1}_{i \in \mathcal{N}} \mathbb{1}_{S_i \neq \emptyset}$  for some parameter  $\varphi^* \in (0, 1)$ .

<sup>25</sup>For example, the prior probability that  $|S| = 1$  could be high enough that mediated learning is infeasible according to Theorem 2, and infinite with very small probability. However, the first witness' signal could reveal that  $|S|$  is infinite, in which case mediated learning could become feasible, as described by Theorem 1.

any round  $i$ , the probability that  $i$  is a witness is given by some function  $\varphi(i, \mathbb{1}_{S_i \neq \emptyset}, m_1^{i-1})$ .

Part (i) guarantees that the more signals have been discovered, the more likely it is that there are no signals left to discover. Part (ii) implies that the identity (witness or investigator) of  $i$  does not reveal information about the amount of evidence left to discover, other than the fact that  $S_i$  was non empty. It rules out situations, for instance, in which the mere fact that  $i$  is a witness implies that there are many other signals left to discover, or the opposite inference in which the fact that  $i$  is *not* a witness implies that there are many more signals to discover.

After observing his signal  $s_i$ , witness  $i$  sends a report  $m_i \in M$ . His realized utility has two parts:

$$U_i(m) = V_i(m) + \epsilon_i(m_i) \quad (4)$$

where  $V_i(m)$  plays the same role as investigators' utility function, and  $\epsilon_i(m_i)$  is a private shock affecting  $i$ 's preferences.

**Assumption 2** *The random variables  $\{\epsilon_i(m_i)\}_{m_i \in M_i}$  are privately observed by  $i$ . Conditional on  $m_1^{i-1}$ , they are independently distributed from one another and from all other exogenous variables in the model. Their density functions  $\{f_{i,m_i}\}_{m_i \in M_i}$  are uniformly bounded above by some arbitrary constant  $\bar{f}$ .*

Assumption 2 allows the possibility that  $i$  have strong public biases, such as  $i$ 's ideology, financial interest or, in a criminal case,  $i$ 's relationship with the defendant. All such public biases are captured by  $V_i$ . Assumption 2 requires that, in addition to any possible public bias,  $i$  is also subject to some (possibly, small) private bias, which may for instance qualify the extent of  $i$ 's public bias.

The bound  $\bar{f}$  appearing in Assumption 2 plays for witnesses the same role as the cost inverse  $1/c$  does for investigators. Intuitively,  $i$ 's disutility of sending a less preferred message  $m_i$  instead of a more preferred message  $m'_i$  is of order  $1/\bar{f}$  in expectation, as explained in Lemma 8.

If  $i$  is a witness, we will say that  $i$ 's message is *uninformative* if it is statistically independent of  $i$ 's signal conditional on  $m_1^{i-1}$ . Otherwise,  $i$ 's message is *informative*. An informative equilibrium is defined as before: there is at least one agent (investigator or witness) who produces an informative message with positive probability. Continuation informative equilibria are defined analogously. The following result is proved in Appendix D.

**Theorem 3** *For any parameters  $(R, c, \lambda)$ , there exist strictly positive thresholds  $\{\underline{F}^k\}_{k \geq 1}$  with the following property:*

*All equilibria are uninformative unless  $F^k \geq \underline{F}^k$  for all  $k \geq 1$ .*

## 2.6 Analysts

Finally, we consider analysts, who differ from investigators along two dimensions:

- Analysts receive a signal for free.
- They incur a (possibly, arbitrarily small) processing cost to learn the content of their signal.

Compared to witnesses, analysts have no biases but they need to make some effort to understand (e.g., analyze) their signal. Compared to investigators, analysts know that they have a signal for sure at hand, so they are not concerned about not finding a signal. However, they still need to incur a cost to learn the content of the signal, like the laboratory testing a blood sample mentioned in the Introduction, which gets the sample at no cost but must process it at some (possibly, small) cost. The model is identical to Section 2.5, except that witnesses are replaced by analysts, as follows: If  $i$  is an analyst, he receives at no cost a signal  $s_i \in S_i$  at the beginning of round  $i$ . To observe the content of  $s_i$ ,  $i$  must incur a positive cost  $c_i$  that is independently distributed from other random variables in the model. The cost  $c_i$  may be arbitrarily small, but we assume that its distribution  $H$  has a bounded density. After learning his cost  $c_i$ ,  $i$  decides whether to observe  $s_i$  (“work”) or to shirk.  $i$  then chooses some message  $m_i$  to send.  $i$ ’s realized utility is given by:

$$U_i(m) = V_i(m) - c_i \mathbb{1}_i \text{ works}$$

where  $V_i(m)$  plays the same role as investigators’ utility function. We maintain Assumption 1. An equilibrium is informative if there exists at least one agent (investigator or analyst) who works with positive probability. The following result is proved in Appendix E.

**Theorem 4** *For any parameters  $(R, c, \lambda)$ , there exist strictly positive thresholds  $\{\underline{F}^k\}_{k \geq 1}$  with the following property:*

*All equilibria are uninformative unless  $F^k \geq \underline{F}^k$  for all  $k \geq 1$ .*

## 3 Ethical Necessity

### 3.1 Ethical Necessity

The arguments presented so far bring us to one of the paper’s key motivating questions: is ethical behavior *necessary* for society to function? Ethical necessity is not a salient concern in economic analysis, which typically focuses on “selfish” agents evolving within the boundary of well-defined institutions, such as market or democratic institutions whose enforcement is taken for granted.<sup>26</sup> In reality, institutions cannot be taken for granted: if agents are unethical, they may attempt to violate these institutions in various ways, and it is unclear why agents should be presumed to act selfishly within the behavioral boundaries imposed by these institutions but ethically with respect to these boundaries. This separation amounts to a “heroic” dichotomy that, at the very least, deserves closer inspection.

By focusing on mediated learning, this paper provides a tractable framework to study ethical necessity.<sup>27</sup> Since mediated learning is required to investigate criminal cases, political corruption, and other possible violations of institutions, successful mediated learning is a necessary condition for society to function: whenever mediated learning requires ethical behavior, so does society.<sup>28</sup>

---

<sup>26</sup>Economists have studied various forms of non-selfish behaviors, such as pro-social behavior arising in dictator and ultimatum games and various forms of altruism. Unlike earlier works, this paper does not study ethical behavior per se, but whether ethical behavior, in the form of an intrinsic motive to seek and reveal the truth, is necessary for society to function. Self-interest remains a central assumption in economic models and is deeply rooted in the discipline. For example, Edgeworth (1881) observed that “self-interest is the first principle of pure economics.”

<sup>27</sup>The present framework relies on a sequential learning structure. One could consider alternative structures. For example, a central agency may ask several intermediaries to seek and report the truth simultaneously and independently of one another. However, the incentives of such a centralized agency must also be scrutinized by a third party, which reintroduces a sequential element. One example concerns potential abuses by various law enforcement agencies, which require the oversight of higher authorities. This and other designs are discussed in Section 5.

<sup>28</sup>In this sense, the paper is related to the work of the philosopher Bernard Williams, who argues that truthfulness must emerge as an intrinsic value in order for the transmission of information to be successful (Williams (2002)). Successful transmission requires that the sender be truthful in the sense of being *accurate* (i.e., actually possess information) and *sincere* (i.e., intend to fully reveal this information) and that the receiver believe that the sender is truthful. When agents have material incentives to lie, Williams expresses the view that agents must value truthfulness intrinsically in order to overcome these material incentives to lie.

While ethical necessity has not played a central role in economic theory, related questions have been considered in previous work.<sup>29</sup> Hurwicz (2007) discusses the existence of “interveners,” defined as ethical monitors in a monitoring-the-monitor problem, and expresses his personal belief in the existence of interveners. Unlike the present paper, however, Hurwicz argues that interveners are not needed for successful monitoring hierarchies.<sup>30</sup> Hurwicz describes an environment with three agents  $A, B, C$ , in which  $B$  monitors  $A$ ’s actions,  $C$  monitors  $B$ ’s monitoring of  $A$ ,  $A$  monitors  $C$ ’s monitoring of  $B$ ’s monitoring  $A$ , and so on. This “Hurwicz triangle,” in which the monitoring hierarchy is folded in a loop going through the agents, omits corruption across monitors, a possibility studied by Strulovici (2021).<sup>31</sup>

Although ethical necessity is not discussed explicitly in economic models of law enforcement, the existence of a reliable monitor is often implicit. In Becker and Stigler’s (1974) study of wrongdoing and malfeasance by enforcement officers, for instance, the authors assume that wrongdoing may be exogenously detected. When ethical monitors are unavailable, however, the question of how this detection is generated remains to be answered.<sup>32</sup>

### 3.2 Existence and Observability of Ethical Agents

The arguments developed so far assume that agents are known to be non-ethical. Without this assumption, mediated learning can be achieved if agents put sufficiently high weight

---

<sup>29</sup>Myerson (2006) shows that in a federalist regime, the existence of virtuous politicians can guarantee the success of democracy either at a national or a provincial level, whereas democracy can fail at both levels when virtuous politicians are surely absent. Glazer and Rubinstein (1998) describe an information aggregation environment in which, if all agents are purely concerned with achieving the social optimum, there always exist equilibria in which the optimum fails to be achieved, whereas if all agents are *also* concerned with their individual recommendation being followed, the social optimum is uniquely selected. Matsushima (2008) studies the possibility of full implementation when agents prefer truth-telling whenever their material payoffs are unaffected by their message.

<sup>30</sup>Rahman (2012) proposes a different approach to monitoring, well suited for repeated monitoring tasks such as controls for airport security: a principal asks agents to sometimes violate the rules on purpose to check whether these violations are caught by the monitor. Such violations are detected at no cost for the principal, since he instigates them. The approach is well suited when violations can be faked at little social cost, the principal has commitment power, and collusion between principal and agents is impossible.

<sup>31</sup>Levine and Modica (2016) consider a similar structure, in which agents in a group take some initial action, then verify with some probability the action previously taken by their neighbor, then verify the earlier verification task of their neighbors, and so forth.

<sup>32</sup>Milgrom, North, and Weingast (1990), who study the enforcement of trade institutions by law merchants in medieval Europe, consider some of the law merchants’ incentives to lie and take bribes.

on the probability that other agents are ethical. To see this, let us consider once more the blood-testing example. If everyone believes that the second agency will behave ethically (i.e., will correctly test the second sample and truthfully reveal the result of its test), this can be used to induce the first agency to truthfully test and report the content of the first blood sample. By comparing the two agencies findings, and punishing the first agency for any discrepancy, the first agency can be held accountable thanks to a second agency's reliable source of information.

More generally, any agent can be incentivized to behave truthfully as long as the agent believes that subsequent agents are likely to provide informative reports.

Conversely, if some agent is ethical but other agents believe that he is not, mediated learning can fail just as when all agents are non-ethical, for two reasons: First, the findings of the ethical agent are (wrongly) believed to be uninformative. Second, precisely because this agent's findings are believed to be uninformative, the findings cannot be used to incentivize other agents to behave truthfully. Therefore, mediated learning requires that agents believe that other agents behave ethically with high enough probability.<sup>33</sup>

The theory thus provides a specific mechanism for why eroding trust in institutions is damaging: society may need everyone to believe in the existence of ethical agents in order to sustain ethical behavior. Events that erode the strength of this belief can have severe consequences for the feasibility of mediated learning and the functioning of society.<sup>34</sup>

Furthermore, it may be empirically difficult to distinguish between agents who have ethical preferences and agents who merely behave ethically because they *believe* in the existence of ethical agents. Put differently, the belief in ethical behavior can be self-fulfilling.<sup>35</sup>

---

<sup>33</sup>It would be interesting to study, following the results of the present paper, how likely and concentrated ethical agents have to be in order to sustain mediated learning.

<sup>34</sup>Belief in ethical behavior gets rid of virtual attrition, but not of real attrition. For example, if the second blood-testing agency is, in fact, non-ethical, it cannot be incentivized to tell the truth: real attrition interrupts mediated learning. But if everyone erroneously believes that the second agency is ethical, then the first agency may be incentivized to behave truthfully.

<sup>35</sup>Bénabou, Falk, and Tirole (2020) study theoretically and empirically the extent to which moral preferences can be elicited.

## 4 Social Consequences

### 4.1 Political Consequences

Consider the perspective of a citizen who is *cynical* in the sense that he does not trust information intermediaries to behave ethically. In this citizen’s mind, agents involved in the learning process do not care about the truth per se and are collectively aware of this.

To see how information attrition affects the views of such a citizen, let us first revisit the blood-testing example. Given that the blood samples are subject to information attrition, it is reasonable—indeed, rational—for a cynical individual to treat as uninformative any finding reported by the blood-testing agencies, for the reasons explained above.<sup>36</sup> As a result, two citizens with cynical beliefs about agencies’ behavior and otherwise different views of the world may rationally entertain very different beliefs about whether a particular athlete used a banned substance, even after agencies have made their reports public: mediated learning fails to convince citizens and to bring their views closer.<sup>37</sup>

Mediated-learning failures can have severe consequences for social cohesion and political stability. When, for instance, a politician is accused of corruption, the average citizen must rely on declarations made by intermediaries, such as officials and journalists, to learn whether the corruption charge is true or, on the contrary, an attempt to smear and neutralize this politician. There is no *deus ex machina* to lift all confusion and reveal the truth to the public. Information attrition is a concern in these environments: incriminating documents can be destroyed, witnesses intimidated or eliminated, and so on.<sup>38</sup>

---

<sup>36</sup>Blood tests could in principle be certified by a third party, which would check whether the agency did its job properly and thus increase the trust in its report. One would then have to understand the third party’s incentives. This “monitoring the monitor” structure is discussed in Section 5.2.

<sup>37</sup>Although the mechanism is unrelated, this result echoes failures of belief convergence in social learning environments. See, e.g., Acemoglu, Chernozhukov, and Woldar (2016).

<sup>38</sup>These observations also apply to government agencies suspected of abusing their power or violating some rules. Examples of evidence destruction by governmental agencies abound, even in prominent democracies. In the United States, for instance, CIA director Richard Helms ordered in 1973 that all documents pertaining to the CIA’s infamous MK Ultra program on mind-control experiments be destroyed (“An interview with Richard Helms”, <https://www.cia.gov/static/9845318ed4b2db36bc185604a2c3bc40/interview-with-richard-helms.pdf>). In 2005, the CIA’s Director of Operations ordered the destruction of all interrogation tapes of Abu Zudaydah and Abd al-Rahim al-Nashiri that featured “enhanced” interrogations (“Tapes by C.I.A. Lived and Died to Save Image,” <https://www.nytimes.com/2007/12/30/washington/30intel.html>.)



Citizens holding different priors about such a corruption case, perhaps due to differing ideologies or political perspectives, but holding a similarly cynical view about the lack of ethical motives underlying public declarations, may therefore both dismiss official statements and journalistic reports about the case and maintain strong disagreements.

The theory thus explains why rational citizens may remain divided on questions subject to information attrition, and why the lack of trust in investigative institutions can rationally perpetuate polarization. It also provides a specific meaning and mechanism for the observation that eroding citizens' "trust" in institutions harms social cohesion.

## 4.2 Historical Revisionism and Information Attrition

Historical revisionism is another instance of mediated learning in which information attrition plays an important role. Understanding historical events is of obvious importance, not only to learn lessons from the past but also to assess claims based on such events, such as territoriality and reparation claims. Mediated learning is necessary because citizens cannot directly verify historical events.<sup>39</sup> They must rely on experts and officials to access and correctly interpret archives, artifacts, and other sources of information. Information attrition is both exogenous (e.g., witnesses die) and endogenous (e.g., documents may be destroyed on purpose).

To give a concrete example,<sup>40</sup> consider the fire of the German parliamentary building (Reichstag) on February 27, 1933. The importance of this event can hardly be overstated. The Nazis, who had lost seats in the previous parliamentary election, claimed that the fire had been caused by communists and used this claim to pressure president Hindenburg into imposing martial law on Germany (the "Reichstag fire decree") and arrest and weaken communists. This allowed the Nazis to form a majority coalition following the March 5, 1933 parliamentary elections, consolidated Hitler's power, and led to the Enabling Act.

While the claim of communist involvement in the fire has long been rejected, there was until recently a consensus among mainstream historians that the Reichstag fire had *not* been

---

<sup>39</sup>This point is obvious with regard to events that took place before citizens' lifetime. Even with regard to contemporaneous events, aggregating information and forming a global picture of an event is a highly complex task that requires expertise, time, and a special access to information. The pitfalls of such information aggregation have been famously illustrated by Stendhal's *La Chartreuse de Parme* whose protagonist, Fabrice del Dongo, takes part to the Battle of Waterloo with the Napoleonic army and construes a completely erroneous version of the battle.

<sup>40</sup>Huq (2018) discusses similar, very recent examples, in which the possibly false threats of terrorism or coups were used to weaken democratic institutions and shift power to more authoritarian regimes.

caused by the Nazis. Fritz Tobias, one of the most respected historians on this subject in the postwar period, published a series of articles purporting to show that van der Lubbe, the person who was convicted for the arson, had acted alone.<sup>41</sup> Historians accepted Tobias's version of the event until 2001, when two historians studying Gestapo archives raised the possibility that it was a group of SA officers who had set the fire (Bahar and Kugel (2001)). Hett (2014) used recent scientific advances to convincingly argue that it would have been impossible for a single individual to set the fire. In 2019, an affidavit written in 1955 by former SA Hans-Martin Lennings was discovered in Tobias's personal files and published by RedaktionsNetzwerk Deutschland, which stated that Lennings and other SAs had driven van der Lubbe from an infirmary to the Reichstag when the fire had already started, effectively setting up van der Lubbe.

Information attrition took several forms: all but one of the SA officers who were allegedly involved in the Reichstag fire were killed (and, hence, silenced) during the Night of the Long Knives; van der Lubbe was beheaded in 1934 for his alleged role in the arson; and any forensic evidence about the Reichstag fire has been long gone. Furthermore, Fritz Tobias hid Lennings' affidavit, which contradicted Tobias' single perpetrator theory, until his death.<sup>42</sup> To this day, the strongest case for Nazi involvement thus seems to come from Lennings' affidavit, and hence hinges on one man's statement. It is unclear what Lennings' motivation for incriminating the SA may have been, except perhaps for setting the record straight. For observers who doubt Lennings' statement, the question of who set the Reichstag on fire may reasonably represent a failure of mediated learning. Lutjens (2016) believes that "the continuous reshaping of the Reichstag fire by those with a stake in the matter has fragmented the truth beyond recovery."

## 5 Discussion

Mediated learning fails if the following conditions hold jointly: (i) intermediaries do not care about the truth, (ii) information is subject to attrition, (iii) there is no exogenous, public revelation of the truth at any future time, and (iv) intermediaries proceed sequentially.

This result holds for arbitrary utility functions as long as intermediaries have arbitrarily

---

<sup>41</sup>The articles were published by *Der Spiegel* under the title „Stehen Sie auf, van der Lubbe!“ in 1959-60.

<sup>42</sup>Tobias is now suspected of protecting former Nazi officers after the war and having a private interest in dissimulating the Nazis' role in the Reichstag fire.

small biases or informational costs. It holds for all equilibria, in the strong sense that nothing at all is learned about the state of the world, and despite the fact that agent incentives can be administered without any further agency problem: whatever rewards and punishments are promised to the agents as a function of reporting histories can be perfectly enforced.

This impossibility result may be viewed as a reference point: to succeed, mediated learning must break at least one of the four conditions above. In particular, mediated learning can succeed if some intermediaries are motivated by the desire to seek the truth, or by a belief that other intermediaries have such a motivation. Several ways out are discussed below.

## 5.1 Escaping Attrition with Hard evidence

In some cases, signals are hard to fabricate and non-disposable. Consider video footage of a crime, in which the criminal is clearly identifiable. Such footage can be viewed numerous times with little degradation and is conceptually similar to reproducible evidence.

Even this kind of evidence is need not be perfectly reliable. For example, video footage can be fabricated, as exemplified by the emergence of baffling deepfakes, which have raised severe concerns (see, e.g., the journalistic work of Schick (2020)).<sup>43</sup>

DNA testing also illustrates this issue. The amount of usable DNA samples on a crime scene is finite. The procedure of DNA testing involves a replication phase, such as PCR amplification or DNA cloning, which creates more “evidence” that can be stored and verified by subsequent investigators. Crucially, however, these replications are only as reliable as the original DNA sample. They do not constitute new, independent evidence.

This is all the more important as DNA samples can be synthesized, i.e., literally “fabricated,” to match any desired DNA profile.<sup>44</sup> Moreover, DNA samples can be erroneously or malevolently taken outside of the crime scene and presented as coming from the scene.<sup>45</sup>

---

<sup>43</sup>Deepfakes can be detected by computer scientists who specialize in the field. However, this kind of detection *increases* the reliance on mediated learning. With deepfakes, video evidence no longer constitutes public, hard evidence. The disappearance of hard video evidence due to deepfakes is emphasized by Schick (2020).

<sup>44</sup>Frumkin, Wasserstrom, Davidson, and Grafit (2010) show the possibility of creating saliva or blood samples with the desired DNA. The authors, as well as subsequent work by other researchers, show that identifying methylation patterns in DNA samples can help distinguish synthetic and natural DNA, although such identification is challenging.

<sup>45</sup>A famous example is the “Phantom of Heilbronn,” a presumed serial killer whose DNA was found on 40

Ultimately, agents joining an investigation can either test the DNA material that previous laboratories have left them, which have been manipulated or fabricated, or look for genuinely new DNA samples, which brings us back to the problem of information attrition.

The role of technology on mediated learning is complex and deserves a separate exploration. While, DNA testing, video footage, and other technological advances have increased the set of reliable evidence, technology can be used to manipulate evidence and do so *more anonymously* than before, *increasing the reliance on experts and, hence, on mediated learning*.

## 5.2 Alternative Learning Structures

Several remedies may be considered to address the unraveling results of Theorems 2, 3, and 4. First, agents could be asked to investigate and report their findings simultaneously, a structure that we will call *parallel monitoring*. Second, agents could investigate the actions of past investigators, rather than the initial fact. Third, while after-the-fact investigations of the type studied in the baseline model are crucial and widely used to maintain ex ante incentives of would-be offenders, in some applications investigators could be incentivized by the perspective that their findings will have an influence on subsequent actions and events, which brings the framework closer to a repeated game setting. Finally, each agent could intervene multiple times in the investigation process, rather than just once. These possibilities are examined in turn.

### 1. Parallel Monitoring with a Centralized Authority

In some applications such as journalistic investigations and academic research, it is realistic to assume that agents report their findings sequentially. It is nonetheless natural to consider, from a mechanism-design perspective, the case in which several agents simultaneously and independently investigate the fact of interest and report their findings to some central authority. A conceptually similar design is to hide past reports from investigators until they have made their own report.<sup>46</sup>

---

crime scenes over a fifteen-year span in Germany, Austria, and France. The DNA turned out to belong to a woman working in the factory that made the cotton swabs used to collect DNA samples.

<sup>46</sup>Gershkov and Szentes (2009) show that such a design is optimal in a voting model with costly information acquisition. In their model, voters' preferences depend on the state of the world and informative equilibria exist even when all reports are public, but optimal learning requires that past reports be hidden.

For example, blood-testing laboratories could be asked to test their respective samples independently from each other and report their findings simultaneously to some overseeing agency. The laboratories would be rewarded if their findings match and punished otherwise.<sup>47</sup>

An important first concern with this solution stems from the incentives of the central authority coordinating the agents. If this authority has a material interest in a specific outcome (which cannot be ruled out, especially in politically charged investigations), it can secretly help agents coordinate their reports or influence agents' reports in various ways. In order to avoid this, the central agency must itself be monitored, which brings us back to a sequential monitoring problem. In the blood-testing example, if laboratories must report their findings simultaneously, there is no recourse for an athlete accused of doping if laboratories conspire to accuse the athlete.

A second issue concerns the robustness of equilibria in which agents are induced to learn. Agents subject to strategic uncertainty may be harder to incentivize. This question is studied by Pei and Strulovici (2022), who find that the answer depends on the concept of robustness used: When one focuses on partial implementation and uses a local, *ex ante* notion of robustness building on Kajii and Morris (1997), it is possible to design robust mechanisms that prompt agents to acquire costly information and report their findings truthfully. For the case of full implementation, or for stronger notions of robustness, however, impossibility results are derived when agents are non-ethical in the sense of the present paper.

In some cases, parallel monitoring may be challenging to implement. It is difficult, for instance, to send multiple investigators on a crime scene to independently interrogate witness and collect evidence, without the investigators being able to communicate, either directly or through witnesses and possibly coordinate. Moreover, when the evidence is in limited supply (such as the weapon of a crime), such a limitation creates negative correlation in investigators' reports, since at most one of them can discover the evidence (weapon).<sup>48</sup>

## 2. Monitoring the Monitor

Instead of all agents investigating the same initial fact, agents could investigate one another. For example, agent 1 would investigate the initial fact, agent 2 would investigate

---

<sup>47</sup>The academic refereeing process has reporting features that resemble the parallel-monitoring design, although the incentives for referees are different and arguably more complex than a mere coordination motive.

<sup>48</sup>The adverse effect of negative correlation on the informativeness of agents with a coordination motive is a central finding in Pei and Strulovici's (2020) analysis of strategic crime.

agent 1’s investigation of the initial fact, agent 3 would investigate agent 2’s investigation of 1, and so on. Such a “monitoring hierarchy” is described by Hurwicz (2007). In his model, the number of agents is finite and the monitoring chain cycles repeatedly across a fixed set of agents. A first conceptual difficulty with this approach concerns the simultaneity and complexity of these monitoring tasks: agents are supposed to conduct an infinite amount of monitoring tasks and are indirectly the subject of the tasks that they are investigating.

Even if we considered an unbounded sequence of distinct agents, each of which is tasked with investigating the previous agent in the sequence, another issue would emerge, which concerns agents’ ability to collude. For example, if an agent discovers incriminating evidence about the agent he was monitoring could hide or destroy the evidence in exchange for a payment from the guilty agent. Such a transfer amounts to a local form of corruption among agents in direct contact. In a separate paper, I show that even such a local form of corruption may suffice to destroy the possibility of mediated learning (Strulovici (2021)).

### 3. Repeated Setting

When mediated learning concerns the identification of a criminal and opportunities to commit crime are repeated over time, it is a priori possible that investigators care about the truth indirectly, through the impact that their findings have on citizens’ future behavior.

Suppose that a citizen’s decision to commit crime depends on whether his past actions were accurately called by past investigators: for example, a citizen who was wrongfully accused in the past or mistakenly acquitted of crimes that he did commit may be more likely to commit crime in the future. In this setting, investigators could in principle have an endogenous incentive to report accurate findings. Provided that players are sufficiently patient, this kind of strategy profile could a priori be used to incentivize accurate reporting.

However, for this argument to work, a citizen’s strategy must depend on his private history, where the public history consists of official findings about citizens’ past actions and a citizen’s private history records his actual past actions. In order for a citizen’s private history to affect his decision of whether to commit crime, the citizen must be indifferent between committing crime and abstaining from it, a knife-edge condition that is violated if, for instance, citizens are subject to small private shocks affecting their benefits from committing crime.<sup>49</sup>

---

<sup>49</sup>The argument is somewhat similar to the section on witnesses in this paper.

## 4. Alternating Statements

Finally, one could ask a fixed set of agents to take turns investigating and reporting on the question of interest. The key difference with the baseline model is that agents now have private information about what they did in the past, which affects how they interpret the declarations of other agents. If an agent has discovered disposable signals in the past, he knows that other agents have fewer signals to discover. The analysis becomes more complex because agents' decisions now depend on their beliefs about the amount of evidence left, about other agents' beliefs about the amount evidence, their belief about agents' beliefs about their beliefs about the amount of evidence left and so on. While information attrition is likely to have a similar effect as in this paper's model, confirming this intuition and exploring this question is left for future research.

## 5.3 Designing Truth-Dependent Preferences: Oaths, Capitalism, and Popular Juries

A more direct approach to improving mediated learning is to increase the salience of truth-dependent preferences.

This may be achieved by fostering agents' ethical sense, from inculcating an ethical education and culture to strengthening trust in institutions and developing effective vetting and selection processes for key learning responsibilities.

Professional oaths, from the Hippocratic oath in medicine to journalistic oaths such as Walter Williams' Journalist's Creed (Farrar (1998)) aim at eliciting ethical behavior. This paper suggests a positive correlation between the need for oaths in various professions and the severity of information attrition in these professions.

Even if a small fraction of agents is swayed by such oaths, this may in principle suffice for incentivizing truthful behavior by other agents. Studying the mechanisms and behavioral features through which ethical agents can incentivize mediated learning is beyond the scope of this paper, but it is easy to conceive of simple examples:<sup>50</sup> suppose that an agent, who is known to truthfully seek and report the truth is commonly known to appear in round  $N > 1$ .

---

<sup>50</sup>There are various approaches to modeling ethical behavior. For instance, Ellingsen and Mohlin (2020) distinguish three dimensions: decency, integrity, and punitivity. Harsanyi (1980) and Feddersen and Sandroni (2006) study rule-utilitarian agents and Roemer (2019) consider Kantian agents.

This agent’s report provides reliable information, akin to an exogenous public signal about the state of the world, which can be used to incentivize all agents coming in rounds  $i < N$ . Even if agent  $N$  has only a small probability  $p < 1$  of behaving ethically, his report may still be used to incentivize agents in earlier rounds as long as these agents’ rewards and punishments are of order  $1/p$ .

Another approach is to organize society in a way that increases information mediators’ material dependence on the truth, i.e., gives them “skin in the game.” Eliciting information from agents about a scientific fact or the social value of a new product or process is easier when the agents stand to gain financially from this information, which may broadly interpreted as capitalistic incentives. Thus interpreted, the theory offers a new perspective on the “virtue” of capitalism relative to systems in which agents have low-powered incentives.<sup>51</sup> The theory also emphasizes that violations of the rules of capitalism (or any other system, for that matter) may be difficult to detect and reveal truthfully, and thus suggests a potential tradeoff between the incentives provided within a given economic or political system and the incentives required to guarantee that the system is respected by its participants.

A final angle to attack mediated learning failures is to democratize the learning process by enlarging the pool of potential information intermediaries. Large pools can increase the alignment—real or perceived—between the intermediaries and society as a whole, in contrast to the baseline model of the paper, in which society’s objective is dissociated from the intermediaries’. Large pools of intermediaries are conceivable when the expertise required to learn the fact of interest is limited. The institution of popular juries may be viewed as one such application, which trades off intermediaries’ expertise with their representativeness of a more global and diffuse body of stakeholders.

---

<sup>51</sup>A large literature emphasizes capitalism’s ability to reduce moral hazard problems relative to socialistic systems. See Myerson (2007) and Tirole (2006) for a review of relevant papers and corporate-finance models capturing this idea. The question of incentive compatibility and its relation to various economic systems is at the heart of Leonid Hurwicz’s development of mechanism design (Hurwicz (1973)).



## References

- ACEMOGLU, D., CHERNOZHUKOV, V., AND YILDIZ, M. (2016) “Fragility of asymptotic agreement under Bayesian learning,” *Theoretical Economics*, Vol. 11, pp. 187–225.
- BAHAR, A. AND KUGEL, W. (2001) *Der Reichstagsbrand. Wie Geschichte gemacht wird*, Edition Q Verlag, Berlin.
- BANERJEE, A. (1992) “A simple model of herd behavior,” *Quarterly Journal of Economics*, Vol. 107, pp. 797–817.
- BARLOW, R., MARSHALL, A., AND PROSCHAN, F. (1963) “Properties of probability distributions with monotone hazard rate,” *Annals of Mathematical Statistics*, Vol. 34, pp. 375–389.
- BECKER, G. AND STIGLER, G. (1974) “Law enforcement, malfeasance, and compensation of enforcers,” *The Journal of Legal Studies*, Vol. 3, pp. 1–18.
- BÉNABOU, R. AND TIROLE, J. (2003) “Intrinsic and extrinsic motivation,” *Review of Economic Studies*, Vol. 70, pp. 489–520.
- BÉNABOU, R., FALK, A., AND TIROLE, J. (2020) “Eliciting moral preferences: Theory and experiment ,” *Working Paper*, Princeton University.
- BIKHCHANDANI, S., HIRSHLEIFER, D. AND WELCH, I. (1992) “A theory of fads, fashion, custom, and cultural change as informational cascades,” *Journal of Political Economy*, Vol. 100, pp. 992–1026.
- EDGEWORTH, F. (1881) *Mathematical psychics: An essay on the application of mathematics to the moral sciences*, Kegan Paul.
- ELLINGSEN, T. AND MOHLIN, E. (2020) “Dutiful behavior: A model of moral sentiments,” *Working Paper*, Stockholm School of Economics.
- FARRAR R. (1998) *A Creed for My Profession: Walter Williams, Journalist to the World*. University of Missouri Press.
- FEDDERSEN, T., SANDRONI, A. (2006) “A theory of participation in elections,” *American Economic Review*, Vol. 96, pp. 1271–1282.
- FRUMKIN, D., WASSERSTROM, A., DAVIDSON, A. AND GRAFIT, A. (2010) “Authentication of forensic DNA samples,” *Forensic science international: genetics*, Vol. 4, pp. 95–103.

- GERSHKOV, A., SZENTES, B. (2009) “Optimal voting schemes with costly information acquisition,” *Journal of Economic Theory*, Vol. 144, pp. 36–68.
- GLAZER, J. AND RUBINSTEIN, A. (1998) “Motives and implementation: On the design of mechanisms to elicit opinions,” *Journal of Economic Theory*, Vol. 79, pp. 157–173.
- GRANOVETTER, M. (2017) *Society and economy*, Harvard University Press.
- GRICE, H.P. (1975) “Logic and conversation,” *In Cole, P., Morgan, J.L. (Eds.), Syntax and Semantics, Vol. 3*, Academic Press, New York, pp. 41–58.
- HARSANYI, J. (1980) “Rule utilitarianism, rights, obligations and the theory of rational behavior,” *Theory and Decision*, Vol. 12, pp. 115–133.
- HETT, B. (2014) *Burning the Reichstag. An Investigation into the Third Reich’s Enduring Mystery*. Oxford University Press.
- HUQ, A. (2018) “Terrorism and democratic recession,” *University of Chicago Law Review*, Vol. 85, pp. 457–484.
- HURWICZ, L. (1973) “The design of mechanisms for resource allocation,” *American Economic Review*, Vol. 63, pp. 1–30.
- HURWICZ, L. (2007) “But who will guard the guardians?” *Nobel Prize Lecture*.
- KAJII, A., AND MORRIS, S. (1997) “The robustness of equilibria to incomplete information,” *Econometrica*, Vol. 65, pp. 1283–1309.
- LAMPORT, L., SHOSTAK, R. AND PEASE, M. (1982) “The Byzantine generals problem,” *ACM Transactions on Programming Languages and Systems*, Vol. 4, pp. 382–401.
- LANDES, W. (1971) “An economic analysis of the courts,” *The Journal of Law and Economics*, Vol. 14, pp. 61–107.
- LEVINE, D., AND MODICA, S. (2016) “Peer discipline and incentives within groups,” *ACM Transactions on Programming Languages and Systems*, Vol. 4, pp. 382–401.
- LUTJENS, R. (2016) “Burning the Reichstag: An investigation into the Third Reich’s enduring mystery by Benjamin Hett (review),” *German Studies Review*, Vol. 39, pp. 411–412.
- MATSUSHIMA, H. (2008) “Role of honesty in full implementation,” *Journal of Economic Theory*, Vol. 139, pp. 353–359.

- MILGROM, P., NORTH, D., AND B. WEINGAST (1990) “The role of institutions in the revival of trade: The law merchant, private judges, and the champagne fairs,” *Economics & Politics*, Vol. 2, pp. 1–23.
- MYERSON, R. (2006) “Federalism and incentives for success of democracy,” *Quarterly Journal of Political Science*, Vol. 1, pp. 3–23.
- MYERSON, R. (2009) “Fundamental theory of institutions: a lecture in honor of Leo Hurwicz,” *Review of Economic Design*, Vol. 13, p. 59–75.
- PEI, H.D., AND STRULOVICI, B. (2020) “Crime entanglement, deterrence, and witness credibility,” *Working Paper*, Northwestern University.
- PEI, H.D., AND STRULOVICI, B. (2022) “Robust implementation with costly information,” *Working Paper*, Northwestern University.
- RAHMAN, D. (2012) “But who will monitor the monitor?,” *American Economic Review*, Vol. 102, pp. 2767–2797.
- ROEMER, J. (2019) *How we cooperate: A theory of Kantian optimization*, Yale University Press.
- SCHICK, N. (2020) *Deepfakes: The coming infocalypse*, Hachette UK.
- SMITH, L., SØRENSEN, P., AND TIAN, J. (2021) “Informational herding, optimal experimentation, and contrarianism” *Review of Economic Studies*, Vol. 88, pp. 2527–2554.
- STRULOVICI, B. (2021) “Learning and corruption on monitoring chains,” *American Economic Association, Papers and Proceedings*, Vol. 111, pp. 544–548.
- TIROLE, J. (2006) *The theory of corporate finance*, Princeton University Press.
- WILLIAMS, B. (2002) *Truth and truthfulness: An essay in genealogy*, Princeton University Press.

# A Proof of Theorem 2

## A.1 Discovery-Belief Formula

The section presents three lemmas leading to Proposition 2, which contains the discovery-belief formula. Lemma 1 states that an agent works only if the probability  $F_i^1$  that there remains at least one signal to discover is high enough. Lemma 2 states that when an agent works and finds nothing, he cannot be much better off, ex post, than if he had shirked instead, especially if  $F_i^1$  is close to 1, which means that the working agent's lack of a discovery is likely attributable to bad luck. Lemma 3 computes an upper bound on an agent's expected benefit from working relative to shirking. This upper bound is expressed in terms of the probability that a working agent produces a message that triggers an informative continuation equilibrium. This set,  $M_i^+$ , plays a key role in the analysis. All three lemmas and Proposition 2 are proved in Appendix C.

We make two simplifications throughout the proof, which are without loss of generality. First, agents' decisions are invariant with respect to a uniform translation in their gross utility functions. We can therefore assume that these functions all take values in some interval  $[0, R]$ . Moreover, since  $R$  is an upper bound on payoffs, it can be increased to guarantee that  $R > c$ , which is assumed from now on.

Let  $\gamma_i$  denote the probability that  $i$  works given  $m_1^{i-1}$ .

**Lemma 1**  $\gamma_i > 0$  only if  $F_i^1 \geq \frac{c}{R} > 0$ .

Given any round  $i$  and report history  $m_1^{i-1}$  such that  $\gamma_i > 0$ , let:

- $V_i^*$  denote  $i$ 's maximal expected gross utility if he shirks, where the maximum is taken over all possible messages  $m_i \in M$  that  $i$  can send after shirking;
- $f_i^0 = 1 - F_i^1$  denote the probability that  $S_i = \emptyset$  at the beginning of round  $i$ .

**Lemma 2** Suppose that  $\lambda < 1$ . If  $i$  works, finds nothing, and sends message  $m_i$ , his expected gross utility  $V_i^w(\emptyset, m_i)$  satisfies

$$V_i^w(\emptyset, m_i) \leq V_i^* + \frac{f_i^0 R}{1 - \lambda}.$$

For any round  $i$  and  $\tilde{M}_i \subset M$ , we consider the following probabilities conditional on history  $m_1^{i-1}$ :

- $g_i(\tilde{M}_i)$ : probability that  $i$  finds a signal and sends a message in  $\tilde{M}_i$  conditional on working;
- $d_i(\tilde{M}_i)$ : probability that  $i$  finds no signal and sends a message in  $\tilde{M}_i$  conditional on working;

Also let  $M_i^+$  denote the set of messages  $m_i$  followed by an informative continuation equilibrium in round  $i + 1$  given the report history  $m_1^{i-1}$  until round  $i$ .

**Lemma 3** *Agent  $i$ 's expected gross utility conditional on working has the following upper bound. If  $\lambda < 1$ , then*

$$V_i^w \leq V_i^* + d_i(M_i^+) \frac{f_i^0 R}{1 - \lambda} + g_i(M_i^+) R.$$

If  $\lambda = 1$ , then

$$V_i^w \leq V_i^* + d_i(M_i^+) R + g_i(M_i^+) R.$$

For each round  $i$ , we introduce the following variables.

- $\beta_i$ : probability that  $i$  discovers a signal conditional on report history  $m_1^{i-1}$  (before observing whether  $i$  works, i.e., viewed from the beginning of round  $i$ );
- $F_{i+1}^k(m_i)$ : probability that there remain at least  $k$  signals at the beginning of round  $i + 1$  given reports  $m_1^i = (m_1^{i-1}, m_i)$ .

**Proposition 2 (DB Formula)** *Let  $C(\lambda) = 2R/(c\lambda(1 - \lambda))$  if  $\lambda < 1$  and  $C(1) = 2R/c$ . The following inequality holds for all constants  $C \geq C(\lambda)$ , round  $i$ , and integer  $k \geq 1$  such that  $F_i^k > CF_i^{k+1}$ :*

$$\beta_i \leq C \frac{\mathbb{E}_i \left[ (F_i^k - F_{i+1}^k(m_i)) \mathbb{1}_{m_i \in M_i^+} \right]}{F_i^k - CF_i^{k+1}}. \quad (5)$$

## A.2 Proof of Theorem 2

The proof proceeds by induction on  $k$ . Lemma 1 already proves the claim for  $k = 1$  with  $\underline{F}^1 = c/R$ . Now suppose that the claim holds for some  $k \geq 1$ : there exists a threshold  $\underline{F}^k > 0$  such that any informative continuation equilibrium in round  $i$  satisfies  $F_i^k \geq \underline{F}^k$ .<sup>52</sup>

<sup>52</sup>Theorem 2, which was stated for round 0, also applies to all continuation equilibria: the thresholds  $\{\underline{F}^k\}_{k \geq 1}$  depend only on the parameters  $(R, c, \lambda)$ , which are constant throughout the game.

We will show that there is a constant  $\underline{F}^{k+1} > 0$  such that a continuation equilibrium can be informative only if  $F_i^{k+1} \geq \underline{F}^{k+1}$ .

For any  $f \in [0, 1]$ , let  $\mathcal{F}^k(f) = \inf\{F_i^k \in [0, 1] : \gamma_i > 0, F_i^{k+1} \leq f\}$ , where the infimum is taken over all on-path histories and all equilibria of the game and is by convention equal to 1 if no informative equilibrium exists for which  $F_i^{k+1} \leq f$ . By construction,  $\mathcal{F}^k(f)$  is nonincreasing in  $f$ .

Let  $\underline{F} = \inf\{f : \mathcal{F}^k(f) < 1\}$ .  $\underline{F}$  is the smallest value of  $F_i^{k+1}$  for which an informative continuation equilibrium exists (or the infimum of such values, if the infimum is not achieved).

Our objective is to show that  $\underline{F} > 0$ . Let  $\bar{F}^k = \lim_{\omega \downarrow \underline{F}} \mathcal{F}^k(\omega)$ , i.e., the right limit of  $\mathcal{F}^k(\cdot)$  at  $\underline{F}$ . This limit is guaranteed to exist because  $\mathcal{F}^k(\cdot)$  is nonincreasing. Intuitively,  $\bar{F}^k$  is defined by looking at all informative equilibria that have the smallest possible probability that there remain at least  $k+1$  signals, and taking the smallest probability that there remain at least  $k$  signals among all such equilibria.

Let  $\varepsilon > 0$  denote any small constant such that  $\hat{F}^k = \mathcal{F}^k(\underline{F} + G\varepsilon) \geq \frac{\bar{F}^k}{1+\eta}$ , where  $G > 0$  is a large constant and  $\eta > 0$  is a small constant, both determined at the end of the proof independently of  $k$ .<sup>53</sup> Since  $\bar{F}^k$  is the right limit of  $\mathcal{F}^k(\cdot)$  at  $\underline{F}$ , such an  $\varepsilon$  exists. Moreover, since all  $\varepsilon' \in (0, \varepsilon)$  also satisfy the condition, we can choose  $\varepsilon$  so that

$$\varepsilon \leq \frac{1}{2} \left( \frac{\hat{F}^k}{G} \right)^2. \quad (6)$$

By definition of  $\underline{F}$ , there exists at least one informative continuation equilibrium for which

$$F_i^{k+1} \in [\underline{F}, \underline{F} + \varepsilon]. \quad (7)$$

Moreover, by definition and monotonicity of  $\mathcal{F}^k(\cdot)$ , there exists an informative continuation equilibrium among those for which  $F_i^k \leq \bar{F}^k + \eta \hat{F}^k$ .

Consider such a continuation equilibrium. Since our objective is to prove that  $\underline{F} > 0$ , suppose by way of contradiction that

$$\underline{F} = 0. \quad (8)$$

From (7), this implies that  $F_i^{k+1} \leq \varepsilon$ . Combining this with Doob's martingale inequality guarantees that  $F_j^{k+1}$  remains small for all  $j > i$  with high probability. Formally, let  $\mathcal{A}$  denote the event that  $F_j^{k+1} \leq G\varepsilon$  for all  $j \geq i$ . We have the following result.

---

<sup>53</sup>For example, if  $\lambda < 1$  one can choose  $G$  and  $\eta$  such that  $\sqrt{G} = 128R^3/(\lambda^2(1-\lambda)c^3)$  and  $\eta = 1/\sqrt{G}$ , as explained at the end of this proof.

**Lemma 4** *i assigns probability at least  $1 - 1/G$  to  $\mathcal{A}$ .*

*Proof.* For  $j \geq i$ , let  $\bar{F}_j^{k+1}$  denote the probability assigned at the beginning of round  $j$  (i.e., conditional on  $m_1^{j-1}$ ) that there remain at least  $k + 1$  signals *at the beginning of round*  $i$  (fixed). The process  $\{\bar{F}_j^{k+1}\}_{j \geq i}$  is nonnegative and bounded above by 1. Moreover, it is a martingale by the law of iterated expectations and the fact that  $j$ 's filtration grows finer as  $j$  increases. Moreover,  $\bar{F}_i^{k+1} = F_i^{k+1} \leq \varepsilon$ .

We can therefore apply Doob's martingale inequality, which implies that for any  $J \geq i$ ,  $\Pr(\max_{i \leq j \leq J} \bar{F}_j^{k+1} \geq G\varepsilon) \leq \frac{\mathbb{E}_i[\bar{F}_i^{k+1}]}{G\varepsilon} \leq 1/G$ . The event  $\bar{\mathcal{A}}_\infty$  defined by  $\{\max_{j \geq i} \bar{F}_j^{k+1} \leq G\varepsilon\}$  is the intersection of the events  $\bar{\mathcal{A}}_J = \{\max_{i \leq j \leq J} \bar{F}_j^{k+1} \leq G\varepsilon\}$ . Therefore,  $\Pr(\bar{\mathcal{A}}_\infty) = \lim_{J \rightarrow \infty} \Pr(\bar{\mathcal{A}}_J) \geq 1 - 1/G$ .

Finally, we note that  $\bar{F}_j^{k+1} \geq F_j^{k+1}$  for all  $j \geq i$ , because the true number of remaining signals only decreases over time and thus whatever signals remained at the beginning of round  $i$  must have contained the signals that remain at the beginning of round  $j$ , so  $\Pr(\mathcal{A}) \geq \Pr(\bar{\mathcal{A}}_\infty) \geq 1 - 1/G$ .  $\blacksquare$

Conditional on  $\mathcal{A}$ ,  $F_j^{k+1} \leq G\varepsilon$  for all  $j \geq i$ . Moreover, if round  $j$  belongs to an informative continuation equilibrium, we have  $F_j^k \geq \mathcal{F}^k(F_j^{k+1}) \geq \mathcal{F}^k(G\varepsilon) = \hat{F}^k$ . Given any positive constant  $C$ , this implies that for informative continuation equilibrium starting in round  $j$ ,

$$F_j^k - CF_j^{k+1} \geq \hat{F}^k - CG\varepsilon \geq \hat{F}^k/2 > 0 \quad (9)$$

provided that  $G \geq C$ , where the last weak inequality comes from (6). The condition  $G \geq C$  will be satisfied by setting  $C = C(\lambda)$ , where  $C(\lambda)$  is defined in Proposition 2, and then choosing  $G$  large enough. (A specific value of  $G$  is given at the end of the proof.)

Let  $\mathcal{A}_j$  denote the event that  $F_l^{k+1} \leq G\varepsilon$  for all integers  $l$  such that  $i \leq l \leq j$ . The events  $\{\mathcal{A}_j\}_{j \geq i}$  form a decreasing sequence that converges to  $\mathcal{A}$  as  $j \rightarrow \infty$ . Moreover,  $\mathcal{A}_j$  is measurable with respect to the information available at the beginning of round  $j$ .

Proposition 2 implies, conditional on event  $\mathcal{A}_j$ , that:

$$\beta_j \leq \mathbb{E}_j \left[ \frac{C(F_j^k - F_{j+1}^k(m_j)) \mathbb{1}_{m_j \in M_j^+}}{F_j^k - CF_j^{k+1}} \right]. \quad (10)$$

Let  $\mathcal{B}_j$  denote the event that  $m_j \in M_j^+$ .  $\mathcal{B}_j$  is adapted to  $m_1^j$ , i.e., known at the beginning of round  $j + 1$ . Notice that if  $\mathcal{B}_j$  does *not* occur, it means by definition that no  $l > j$  ever works (i.e.,  $m_1^j$  is followed by an *uninformative* continuation equilibrium). This implies that the

sequence of events  $\{\mathcal{B}_j\}_{j \geq i}$  is decreasing path by path and, hence, that the sequence  $\{\mathbb{1}_{\mathcal{B}_j}\}_{j \geq i}$  is nonincreasing.

Since  $\mathcal{A}_j$  is measurable with respect to  $m_1^{j-1}$ , Equations (9) and (10) imply that for  $j \geq i+1$

$$\mathbb{E}_{i+1}[\mathbb{1}_{\mathcal{A}_j} \beta_j] \leq \mathbb{E}_{i+1} \left[ \mathbb{E}_j \left[ \mathbb{1}_{\mathcal{A}_j} \frac{2C(F_j^k - F_{j+1}^k) \mathbb{1}_{\mathcal{B}_j}}{\hat{F}^k} \right] \right].$$

The law of iterated expectations then implies that

$$\mathbb{E}_{i+1}[\mathbb{1}_{\mathcal{A}_j} \beta_j] \leq \mathbb{E}_{i+1} \left[ \mathbb{1}_{\mathcal{A}_j} \frac{2C(F_j^k - F_{j+1}^k) \mathbb{1}_{\mathcal{B}_j}}{\hat{F}^k} \right] = \frac{2C}{\hat{F}^k} \mathbb{E}_{i+1}[\mathbb{1}_{\mathcal{A}_j} (F_j^k - F_{j+1}^k) \mathbb{1}_{\mathcal{B}_j}]. \quad (11)$$

Summing (11) over all  $j \geq i+1$ , we obtain

$$\mathbb{E}_{i+1} \left[ \sum_{j \geq i+1} \mathbb{1}_{\mathcal{A}_j} \beta_j \right] \leq \frac{2C}{\hat{F}^k} \mathbb{E}_{i+1} \left[ \sum_{j \geq i+1} \mathbb{1}_{\mathcal{A}_j} \mathbb{1}_{\mathcal{B}_j} (F_j^k - F_{j+1}^k(m_j)) \right]. \quad (12)$$

Let  $J$  denote the first (possibly infinite) round for which the product of the indicator functions in the right-hand side of (12) is zero:

$$J = \inf\{j \geq i+1 : \mathbb{1}_{\mathcal{A}_j} \mathbb{1}_{\mathcal{B}_j} = 0\}$$

with the convention that  $J = +\infty$  if the set is empty. Since both of these indicator functions are nonincreasing in  $j$  for  $j \geq i+1$ , path by path, their product must be equal to zero for all  $j \geq J$ . In words,  $J$  is either the first round in which  $F_j^{k+1}$  exceeds  $G\varepsilon$ , or the *last* round at which the continuation equilibrium is informative (which implies that  $\gamma_j = 0$  for all  $j > J$ ), whichever occurs first.<sup>54</sup>

Consider first the paths for which  $J$  is finite. The argument of  $\mathbb{E}_{i+1}[\cdot]$  on the right-hand side of (12) then reduces to

$$\sum_{j=i+1}^{J-1} (F_j^k - F_{j+1}^k) = F_{i+1}^k - F_J^k.$$

Consider now paths for which  $J$  is infinite. In this case, the argument of  $E_{i+1}[\cdot]$  on the right-hand side of (12) is equal to

$$\sum_{j=i+1}^{\infty} (F_j^k - F_{j+1}^k) = \lim_{\tilde{J} \rightarrow \infty} \sum_{j=i+1}^{\tilde{J}} (F_j^k - F_{j+1}^k) = \lim_{\tilde{J} \rightarrow \infty} \{F_{i+1}^k - F_{\tilde{J}+1}^k\} = F_{i+1}^k - \lim_{\tilde{J} \rightarrow \infty} F_{\tilde{J}}^k.$$

---

<sup>54</sup> $J$  is not a stopping time with respect to the filtration  $\{\mathbb{F}_j\}_{j \geq i+1}$  generated by public histories  $\{m_1^{j-1}\}_{j \geq i+1}$ , because at the beginning of any round  $j$  at which  $j$  works with positive probability, it is unknown whether  $j$  will be the last round in which the agent works. For this reason, the proofs below do not rely on any theorems pertaining to stopping times, such as the optional sampling theorem or the strong Markov property.



Notice that the limit is well defined because  $F_j^k$  is a nonnegative supermartingale with respect to  $j$ .<sup>55</sup> We will call this limit  $F_j^k$  to be consistent with the case in which  $J$  is finite.

Combining these observations with (12), we conclude that

$$\mathbb{E}_{i+1} \left[ \sum_{j>i} \mathbb{1}_{\mathcal{A}_j} \beta_j \right] \leq \frac{2C}{\hat{F}^k} (F_{i+1}^k - \mathbb{E}_{i+1}[F_J^k]).$$

Since  $\beta_j$  is the probability that  $j$  discovers a signal conditionally on  $m_1^{i-1}$ , it may be decomposed as  $\beta_j = \gamma_j F_j^1 \lambda$ , because  $j$  discovers a signal only if three independent (conditionally on  $m_1^{i-1}$ ) events occur: (i)  $j$  works, (ii)  $S_i$  is non empty, and (iii)  $j$  is “lucky” to discover such an element from  $S_i$ . Moreover, Lemma 1 shows that  $\gamma_j > 0$  only if  $F_j^1 \geq c/R$ . This implies<sup>56</sup> that  $\gamma_j \leq g\beta_j$  where  $g = \frac{R}{\lambda c}$  and hence that

$$\mathbb{E}_{i+1} \left[ \sum_{j \geq i+1} \mathbb{1}_{\mathcal{A}_j} \gamma_j \right] \leq \frac{2Cg}{\hat{F}^k} (F_{i+1}^k - \mathbb{E}_{i+1}[F_J^k]). \quad (13)$$

Let  $\mathcal{Z}$  denote the event that at least one agent  $j > i$  works and  $\pi_{i+1}(m_i) = \Pr_{i+1}(\mathcal{A} \cap \mathcal{Z})$ ,

---

<sup>55</sup>The argument is similar to the one used to prove Lemma 4. For any fixed  $j$ , let  $\bar{F}_l^{k+1}$  denote the probability assigned by  $l \geq j$  to there remaining at least  $k+1$  signals *at the beginning of round  $j$* . The process  $\{\bar{F}_l^{k+1}\}_{l \geq j}$  is a martingale in  $j$ , by the law of iterated expectations and the fact that  $l$ 's filtration grows finer as  $l$  increases. Moreover,  $\bar{F}_l^{k+1} \geq F_l^{k+1}$  for all  $l \geq j$ , because the actual number of remaining signals only decreases over time. Therefore, we have  $F_j^{k+1} = \bar{F}_j^{k+1} = E_j[\bar{F}_{j+1}^{k+1}] \geq E_j[F_{j+1}^{k+1}]$ . This, together with the fact that  $F_l^{k+1}$  is uniformly bounded and measurable with respect to the information at the beginning round  $l$ , shows that it is a supermartingale.

<sup>56</sup>The inequality clearly holds if  $\gamma_j = 0$ .

i.e., the probability that  $\mathcal{A}$  and  $\mathcal{Z}$  both occur conditional on  $m_1^i$ . We have

$$\begin{aligned}
\pi_{i+1}(m_i) &= \Pr_{i+1} \left( \mathbb{1}_{\mathcal{A}} \sum_{j>i} \mathbb{1}_{j \text{ works}} \geq 1 \right) \\
&\leq \mathbb{E}_{i+1} \left[ \mathbb{1}_{\mathcal{A}} \sum_{j>i} \mathbb{1}_{j \text{ works}} \right] \\
&= \sum_{j>i} \mathbb{E}_{i+1} \left[ \mathbb{1}_{\mathcal{A}} \mathbb{1}_{j \text{ works}} \right] \\
&\leq \sum_{j>i} \mathbb{E}_{i+1} \left[ \mathbb{1}_{\mathcal{A}_j} \mathbb{1}_{j \text{ works}} \right] \\
&= \sum_{j>i} \mathbb{E}_{i+1} \left[ \mathbb{E}_j \left[ \mathbb{1}_{\mathcal{A}_j} \mathbb{1}_{j \text{ works}} \right] \right] \\
&= \sum_{j>i} \mathbb{E}_{i+1} \left[ \mathbb{1}_{\mathcal{A}_j} \mathbb{E}_j \left[ \mathbb{1}_{j \text{ works}} \right] \right] \\
&= \sum_{j>i} \mathbb{E}_{i+1} \left[ \mathbb{1}_{\mathcal{A}_j} \gamma_j \right].
\end{aligned}$$

The first equality comes from the fact that  $\mathcal{Z}$  is identical to the event  $\{\sum_{j>i} \mathbb{1}_{j \text{ works}} \geq 1\}$ . The first inequality comes from the fact that the random variable  $\mathbb{1}_{\mathcal{A}} \sum_{j>i} \mathbb{1}_{j \text{ works}}$  is nonnegative and integer valued, which implies that its expectation exceeds the probability that it is strictly positive. The second equality is an application of Tonelli's theorem. The second inequality comes from the fact that  $\mathcal{A} \subset \mathcal{A}_j$  and, hence,  $\mathbb{1}_{\mathcal{A}} \leq \mathbb{1}_{\mathcal{A}_j}$ . The next equality comes from the law of iterated expectations and the next one comes from the fact that  $\mathcal{A}_j$  is measurable with respect to the information at the beginning of round  $j$ . The last equality holds by definition of  $\gamma_j$ .

From (13) and Tonelli's theorem, this implies that

$$\pi_{i+1}(m_i) \leq \frac{2Cg}{\hat{F}^k} (F_{i+1}^k - \mathbb{E}_{i+1}[F_j^k]). \quad (14)$$

Let  $F_i^{k,r}(m_i)$  denote the probability that there are at least  $k$  signals left conditional on  $i$  working and reporting  $m_i$ . For any signal  $s_i$ , let  $F_i^{k,w}(s_i)$  denote the probability that there are at least  $k$  signals left conditional on  $i$  working and discovering  $s_i$ .  $F_i^{k,w}(s_i)$  represents  $i$ 's belief after discovering  $s_i$ , while  $F_i^{k,r}(m_i)$  represents what  $i+1$  would believe after observing report  $m_i$  if he knew for sure that  $i$  has worked.

The following lemmas are proved in Appendices C.5 and C.6. Let  $N_i = \{m_i : F_i^{k,r}(m_i) > (1 + \eta)F_i^k\}$ .

**Lemma 5** (i)  $\gamma_i(N_i) \leq \frac{1}{2\eta G^2}$ , (ii) for  $m_i \notin N_i$ ,  $\pi_{i+1}(m_i) \leq \frac{2Cg}{\hat{F}^k} ((1 + \eta)F_i^k - \mathbb{E}_{i+1}[F_J^k])$ .

Let  $T_i = \{m_i : F_{i+1}^{k+1}(m_i) \geq \sqrt{G}\varepsilon\}$ .

**Lemma 6** (i)  $\Pr_{i+1}(\mathcal{A}) \geq 1 - 1/\sqrt{G}$  for all  $m_i \notin T_i$ , (ii)  $\gamma_i(T_i) \leq 1/\sqrt{G}$ .

Let  $V^w$  denote  $i$ 's expected gross utility if he works and  $V^w(m_i)$  denote his expected gross utility conditional on working and reporting  $m_i$ . We have

$$\begin{aligned} V^w &= \sum_{i \in M_i} \gamma_i(m_i) V^w(m_i) \\ &\leq (\gamma_i(N_i) + \gamma_i(T_i))R + \sum_{m_i \notin N_i \cup T_i} \gamma_i(m_i) V^w(m_i) \\ &\leq \left( \frac{1}{2\eta G^2} + \frac{1}{\sqrt{G}} \right) R + \sum_{m_i \notin N_i \cup T_i} \gamma_i(m_i) V^w(m_i) \end{aligned} \quad (15)$$

Moreover,

$$V^w(m_i) = \Pr(\mathcal{Z}|m_1^i) V^w(m_i|\mathcal{Z}) + \Pr(\mathcal{Z}^c|m_1^i) V^w(m_i|\mathcal{Z}^c). \quad (16)$$

Conditional on  $m_1^i$ , the event  $\mathcal{Z}$  is independent of how  $m_i$  was produced (i.e., whether  $m_i$  was obtained by work or fabrication). Indeed, as long as no one works, the distribution of reports  $m_j$  made by agents following  $i$  depends only on  $m_1^i$ , not on the signals that remain to be discovered in the case. And as soon as someone works, then by definition  $\mathcal{Z}$  has occurred. Thus, what triggers the event  $\mathcal{Z}$  (whenever it occurs) is a sequence of uninformative (until  $\mathcal{Z}$  occurs) reports  $m_j$  for agents following  $i$ , whose probability distribution is completely pinned down by  $m_1^i$ .

From the previous lemmas, we have  $\Pr(\mathcal{A} \cap \mathcal{Z}|m_1^i) \leq \frac{2Cg}{\hat{F}^k} ((1 + \eta)F_i^k - \mathbb{E}_{i+1}[F_J^k])$  for all  $m_i \notin N_i$  and  $\Pr(\mathcal{A}^c|m_1^i) \leq 1/\sqrt{G}$  for all  $m_i \notin T_i$ . Letting  $\hat{M}_i = M_i \setminus (N_i \cup T_i)$ , this implies that

$$\Pr(\mathcal{Z}|m_1^i) = \Pr(\mathcal{Z} \cap \mathcal{A}|m_1^i) + \Pr(\mathcal{Z} \cap \mathcal{A}^c|m_1^i) \leq \frac{2Cg}{\hat{F}^k} (F_i^k(1 + \eta) - \mathbb{E}_{i+1}[F_J^k]) + 1/\sqrt{G} \quad (17)$$

for all  $m_i \in \hat{M}_i$ .

Conditional on  $\mathcal{A}$ ,  $F_j^{k+1} \leq G\varepsilon$  for all  $j \geq i$ . By definition of  $J$ , all continuation equilibria until round  $J$  included are informative, which implies that  $F_j^k \geq \mathcal{F}^k(F_j^{k+1})$  for all  $j \leq J$ . Since  $\mathcal{F}^k(\cdot)$  is nonincreasing, this implies that  $F_j^k \geq \mathcal{F}^k(G\varepsilon) = \hat{F}^k$  for all  $j \leq J$ .

We thus have for  $m_i \in \hat{M}_i$

$$\begin{aligned}\mathbb{E}_{i+1}F_j^k &= \Pr_{i+1}(\mathcal{A}) \mathbb{E}_{i+1}[F_j^k|\mathcal{A}] + \Pr_{i+1}(\mathcal{A}^c) \mathbb{E}_{i+1}[F_j^k|\mathcal{A}^c] \\ &\geq \Pr_{i+1}(\mathcal{A})\mathbb{E}_{i+1}[F_j^k|\mathcal{A}] \\ &\geq (1 - 1/\sqrt{G})\hat{F}^k.\end{aligned}$$

By construction,  $F_i^k \leq \bar{F}^k + \hat{F}^k\eta$  and  $\hat{F}^k \geq \bar{F}^k - \hat{F}^k\eta$ . Therefore,  $F_i^k - \hat{F}^k \leq (\bar{F}^k + \hat{F}^k\eta) - (\bar{F}^k - \hat{F}^k\eta) = 2\eta\hat{F}^k$ . Letting  $B = 8Cg$ , (17) then implies (for  $\eta \leq 1$ , which we assume) that for all  $m_i$  in  $\hat{M}_i$

$$\Pr(\mathcal{Z}|m_1^i) \leq B\eta + \frac{B}{2\sqrt{G}} + 1/\sqrt{G}. \quad (18)$$

For each  $m_i$ , let  $V_i^f(m_i|\mathcal{Z}^c)$  denote  $i$ 's expected gross utility if he sends message  $m_i$  conditional on no  $j > i$  working and  $V^{f,*}$  denote the maximizer of  $V_i^f(m_i|\mathcal{Z}^c)$  over all messages  $m_i \in \hat{M}_i$ . Notice that  $i$ 's expected gross utility conditional on  $m_i$  and no  $j > i$  working does not depend on whether  $i$  worked or shirked: either way, the subsequent reports  $\{m_j\}_{j>i}$  are independent of the signals that remain to be discovered. Therefore,  $i$ 's conditional expected gross utilities satisfy  $V^w(m_i|\mathcal{Z}^c) = V_i^f(m_i|\mathcal{Z}^c)$ .

Combining these observations with (16) and (18), we obtain

$$V^w(m_i) \leq \left( B\eta + \frac{B}{2\sqrt{G}} + \frac{1}{\sqrt{G}} \right) R + V^{f,*},$$

for  $m_i \in \hat{M}_i$ . Combining this with (15) yields

$$V^w \leq \left( \frac{1}{2\eta G^2} + \frac{1}{\sqrt{G}} \right) R + \left( B\eta + \frac{B}{2\sqrt{G}} + \frac{1}{\sqrt{G}} \right) R + V^{f,*}.$$

$i$ 's utility from working thus satisfies

$$U^w \leq \left( \frac{1}{2\eta G^2} + \frac{1}{\sqrt{G}} + B\eta + \frac{B}{2\sqrt{G}} + \frac{1}{\sqrt{G}} \right) R + V^{f,*} - c. \quad (19)$$

If  $i$  sends a message  $m_i^* \in \hat{M}_i$  that achieves  $V^{f,*}$ , his utility  $U^f$  satisfies

$$\begin{aligned}U^f &\geq \Pr(\mathcal{Z}|m_1^{i-1}, m_i^*) \times 0 + \Pr(\mathcal{Z}^c|m_1^{i-1}, m_i^*)V^{f,*} \\ &\geq \left( 1 - B\eta - \frac{B}{2\sqrt{G}} - \frac{1}{\sqrt{G}} \right) V^{f,*} \\ &\geq V^{f,*} - \left( B\eta - \frac{B}{2\sqrt{G}} - \frac{1}{\sqrt{G}} \right) R\end{aligned}$$

where 0 is used as a lower bound on  $i$ 's realized gross utility in the first inequality.

Therefore, working is strictly suboptimal if

$$\left( \frac{1}{2\eta G^2} + \frac{1}{\sqrt{G}} + B\eta + \frac{B}{2\sqrt{G}} + \frac{1}{\sqrt{G}} \right) R + V^{f,*} - c < V^{f*} - \left( B\eta + \frac{B}{2\sqrt{G}} + \frac{1}{\sqrt{G}} \right) R$$

or

$$\frac{c}{R} > \frac{1}{2\eta G^2} + \frac{3}{\sqrt{G}} + 2B\eta + \frac{B}{\sqrt{G}}. \quad (20)$$

This inequality is always satisfied as long as  $\eta$  is small enough and  $\eta G$  is large enough. In particular, since  $B = 8Cg$ ,  $g = R/\lambda c$ , and  $C$  can be taken to equal  $2R/c$  if  $\lambda = 1$  and  $2R/(c\lambda(1 - \lambda))$  as noted in Proposition 2, the inequality is satisfied if  $\eta = 1/\sqrt{G}$  and  $\sqrt{G} = 128R^3/c^3$  for  $\lambda = 1$  and  $\sqrt{G} = 128R^3/(c^3\lambda^2(1 - \lambda))$  if  $\lambda < 1$  (recalling that  $R > c$ ), in which case each term on the right-hand side of (20) is less than  $c/4R$  with some strict inequalities. ■

# ONLINE APPENDIX

## B Proof of Theorem 1

To prove Theorem 1, we need to find, for any given  $\rho \in (0, 1]$  and  $\lambda \in (0, 1]$ , some joint distribution of  $(\omega, S)$  that satisfies the condition of the theorem and some payoff functions for which an informative equilibrium exist.

We first prove the theorem when  $\lambda = 1$  and  $\rho = 1$ , which means that any agent who works finds a signal with probability 1. For this case, it suffices to prove Proposition 1.

### B.1 Proof of Proposition 1

We construct payoff functions for which the strategy profile described in the proposition constitutes an equilibrium. Under this strategy profile, as long as  $p_i$  lies in  $(\underline{p}, \bar{p})$  all agents work and truthfully report their signal about  $\omega$ . Moreover, given the symmetric signal structure,  $p_i$  depends only on the number of “H” and “L” signals as long as all agents  $j < i$  work with probability 1. Therefore, the set of equilibrium posteriors forms a grid  $\{q^k\}$  containing  $\hat{p}$  and containing a single point on each side of  $(\underline{p}, \bar{p})$ . Let  $q^0 \leq \underline{p} < q^1, \dots, \hat{p}, \dots, q^N < \bar{p} < q^{N+1}$  denote this grid. Along the candidate equilibrium, the belief  $p_i$  evolves on this grid until it hits either  $q^0$  or  $q^{N+1}$ , after which the investigation stops.

Let  $J$  denote the last investigator who works: we have  $p_J \in \{q^1, q^N\}$  and  $p_{J+1} \in \{q^0, q^{N+1}\}$ . Also let  $\tilde{p} = p_{J+1}$  denote the value of the belief when learning stops under the candidate equilibrium.

We construct utility functions in which an investigator’s compensation depends only on his report and on the posterior  $\tilde{p}$ .

For any  $i$  such that  $p_i = q^k \in (\underline{p}, \bar{p})$ , if  $i$  reports “H”, he receives a reward  $R_H^k \geq 0$  if  $\tilde{p} = q^{N+1}$  and a punishment  $P_L^k \leq 0$  if  $\tilde{p} = q^0$ . If  $i$  reports “L”, he gets  $R_L^k \geq 0$  if  $\tilde{p} = q^0$  and  $P_L^k \leq 0$  if  $\tilde{p} = q^{N+1}$ .

For any  $p, q$  on the grid, let  $\pi(p, q)$  denote the probability that the belief sequence ends with  $\tilde{p} = q^{N+1}$ , i.e., exits  $(\underline{p}, \bar{p})$  through  $\bar{p}$ , from the perspective of an agent who assigns probability  $p$  to  $\omega$  but when the prior used by investigators is  $p_0 = q$ . That is,  $\pi(p, q)$  is the probability

that an individual with prior  $p$  assigns to the sequence  $p_i$  converging to  $q^{N+1}$  in equilibrium given that the public belief, which serves as the state variable for the equilibrium, starts with prior belief  $q$ .

If  $i$  sends report “ $H$ ” starting from prior  $p_i = q^k$ , he assigns a probability  $\pi(q^k, q^{k+1})$  to the public belief converging to  $q^{N+1}$ . If  $i$  works and receives report “ $H$ ”, his belief about the continuation equilibrium is  $\pi(q^{k+1}, q^{k+1})$ . Similarly, if  $i$  sends “ $L$ ”, his belief is  $\pi(q^k, q^{k-1})$  whereas if he works and reports “ $L$ ” his belief is  $\pi(q^{k-1}, q^{k-1})$ . It is straightforward to verify the inequalities

$$\pi(q^{k+1}, q^{k+1}) > \pi(q^k, q^{k+1}) \quad (21)$$

and

$$\pi(q^{k-1}, q^{k-1}) < \pi(q^k, q^{k-1}), \quad (22)$$

for all  $k \in [2, N - 1]$ . The strictness of the inequalities comes from the fact that conditional on the true state  $\omega$ , the path of  $\{p_j\}_{j \geq i+1}$  starting any given value of  $p_{i+1}$  is strictly increasing in  $\omega$  in FOSD, as is easily checked.<sup>57</sup> Therefore, the probability of hitting  $q^{N+1}$  before  $q^0$  is strictly increasing in the belief  $p_i$  that the state is high.

For  $k = 1$ , the investigation stops if  $i$  reports “ $L$ ” so (22) holds as an equality, but (21) is still strict, because this report triggers further investigation. The reverse is true for  $k = N$ : (21) only holds as an equality while (22) is strict.

If  $i$  shirks, his maximal utility is

$$\max\{\pi(q^k, q^{k+1})R_H^k + (1 - \pi(q^k, q^{k+1}))P_H^k; \pi(q^k, q^{k-1})P_L^k + (1 - \pi(q^k, q^{k-1}))R_L^k\}. \quad (23)$$

The left argument is  $i$ 's expected payoff if he sends “ $H$ ”, and the right one is his payoff if he sends “ $L$ ”. Since  $i$  can send either message at no cost, his highest payoff from fabrication is the maximum of these two terms. If  $i$  works, he gets

$$z^k[\pi(q^{k+1}, q^{k+1})R_H^k + (1 - \pi(q^k, q^{k+1}))P_H^k] + (1 - z^k)[\pi(q^{k-1}, q^{k-1})P_L^k + (1 - \pi(q^{k-1}, q^{k-1}))R_L^k] \quad (24)$$

where  $z^k$  is the probability of receiving signal “ $H$ ” given belief  $q^k$ , and is equal to  $z^k = \Pr(“H”|q^k) = q^k\pi + (1 - q^k) \times (1 - \pi)$ .

Working is optimal for  $i$  if (24) exceeds (23) by at least  $c$ .

---

<sup>57</sup>There is a probability space in which the following property holds: if a working agent gets a high signal when  $\omega = L$ , then he must also get a high signal when  $\omega = H$ . In this probability space, the sequence  $\{p_j\}_{j \geq i+1}$  is pointwise higher for  $\omega = H$  than for  $\omega = L$ .

This condition is obtained as follows: set  $P_H^k = P_L^k = -Q$  where  $Q$  is a strictly positive constant, and let  $R_H^k = Q \frac{1-\pi(q^k, q^{k+1})}{\pi(q^k, q^{k+1})}$  and  $R_L^k = Q \frac{\pi(q^k, q^{k-1})}{1-\pi(q^k, q^{k-1})}$ . This guarantees that  $i$ 's expected payoff from fabrication is zero, regardless of the outcome. From (21) and (22), his payoff from working is of order  $Q$  and thus exceeds  $c$ , for  $Q$  high enough.<sup>58</sup>

If  $k = 1$  or  $N$ , there is one signal that  $i$  can send after working which yields a payoff of order  $Q$ , while the other signal yields 0. The signal associated with a positive payoff arises with a probability that is bounded away from 0, since  $p_i$  lies in  $(\underline{p}, \bar{p})$ .

Moreover this scheme is feasible as long as the maximal reward  $R$  and punishment  $-R$  respectively exceed  $\sup\{R_\theta^k : \theta \in \{L, H\}, k \in \{1, \dots, N\}\}$  and  $Q$ .

## B.2 Proof of Theorem 1: General case

The argument of Section B extends easily when  $\rho$  and/or  $\lambda$  are less than 1.

With  $\rho < 1$  and  $\lambda = 1$ , the informative equilibrium is identical to the one described in Proposition 1 except that learning stops as soon as an agent fails to report evidence, in which case he gets a zero compensation. By construction of the equilibrium in Section B.1, shirking and reporting that no evidence was found has the same value as fabricating any other message and can thus be deterred. Since a working agent may find nothing, or the learning process may be interrupted before the belief process exits  $(\underline{p}, \bar{p})$ , in which case the working agent receives 0, the rewards and punishments must be scaled up by  $1/\pi_\rho(q^k)$ , where  $\pi_\rho(q^k)$  is the probability that the belief process exits  $(\underline{p}, \bar{p})$  in equilibrium, given the current belief  $q^k$ , so that the expected compensation of a working agent still exceeds the cost  $c$  of working.

If  $\lambda < 1$  and  $\rho = 1$ , a working agent may fail to find evidence even when there surely exists some. In this case, we assume once more that the compensation is zero, which deters shirking and reporting the empty message, and we scale up all rewards and punishments by  $1/\lambda$  to incentive the agent to work, as in the previous paragraph. The belief process will surely exit  $(\underline{p}, \bar{p})$  since the amount of evidence is unlimited (only individual agents may be unlucky and find nothing with probability  $1 - \lambda$ ).

The case in which both  $\lambda$  and  $\rho$  are less than 1 is a convex combination of the previous cases

---

<sup>58</sup>To see this, let  $\bar{\pi}$  denote a strictly positive lower bound on all inequalities (21) and (22) over all  $k$ 's whenever they hold strictly. Then, the gain from working is of order  $Q\bar{\pi}$ .



and addressed accordingly.

## C Remaining Proofs for Theorem 2

### C.1 Proof of Lemma 1

Let  $V_i^w$  denote  $i$ 's expected gross utility if he works.  $V_i^w$  may be decomposed in terms of  $i$ 's expected gross utility  $\bar{V}_i$  if he works *and* there exists some signal left to be found, and his expected gross utility  $V_i^{w,\emptyset}$  if he works but there is no signal left to be found ( $S_i = \emptyset$ ):

$$V_i^w = F_i^1 \bar{V}_i + (1 - F_i^1) V_i^{w,\emptyset}.$$

Conditional on  $S_i = \emptyset$ ,  $i$ 's expected gross utility if he works is the same as his expected gross utility  $V_i^{s,\emptyset}$  if he shirks and uses the same reporting strategy as he does after working and finding nothing: conditional on  $i$ 's report (whatever it is), the distribution of reports by subsequent agents is identical since there is no signal left to be found. Therefore,  $V_i^{w,\emptyset} = V_i^{s,\emptyset}$ . Furthermore, we also have  $\bar{V}_i \leq R$  since  $R$  is the maximum possible gross utility.

Therefore,  $i$ 's net utility  $U_i^w$  from working, including the cost of working, satisfies

$$U_i^w \leq F_i^1 R + (1 - F_i^1) V_i^{s,\emptyset} - c.$$

Similarly,  $i$ 's utility  $U_i^s$  from shirking satisfies

$$U_i^s \geq F_i^1 \times 0 + (1 - F_i^1) V_i^{s,\emptyset} = (1 - F_i^1) V_i^{s,\emptyset}$$

where the inequality comes from the fact that 0 is a lower bound on  $i$ 's realized gross utility. Comparing the previous two equations shows that shirking strictly dominates working if  $F_i^1 R - c < 0$ .

### C.2 Proof of Lemma 2

For any sequence  $S''$  of signals, let  $\Delta_i(S'')$  denote the probability that  $S_i = S''$  conditional on report history  $m_1^{i-1}$ , and  $\Delta_i^\emptyset(S'')$  denote the probability that  $S_i = S''$  conditional on  $i$  working and finding nothing. Bayesian updating implies that for any  $S'' \neq \emptyset$ :

$$\Delta_i^\emptyset(S'') = \Delta_i(S'') \frac{(1 - \lambda)}{(1 - f_i^0)(1 - \lambda) + f_i^0},$$

and for  $S'' = \emptyset$ :

$$\Delta_i^\emptyset(\emptyset) = \Delta_i(\emptyset) \frac{1}{(1 - f_i^0)(1 - \lambda) + f_i^0}.$$

This implies that

$$\Delta_i^\emptyset(S'') - \Delta_i(S'') = -\frac{\lambda f_i^0 \Delta_i(S'')}{(1 - f_i^0)(1 - \lambda) + f_i^0} \quad (25)$$

for any  $S'' \neq \emptyset$ , and

$$\Delta_i^\emptyset(\emptyset) - \Delta_i(\emptyset) = \frac{\lambda(1 - f_i^0)\Delta_i(\emptyset)}{(1 - f_i^0)(1 - \lambda) + f_i^0}. \quad (26)$$

Let  $V_i(m_i, S'')$  denote  $i$ 's expected gross utility conditional on  $i$  producing evidence  $m_i$  and on  $S_{i+1} = S''$ . Notice that  $m_1^i = (m_1^{i-1}, m_i)$  and  $S_{i+1}$  completely determine the distribution of reports  $\{m_j\}_{j>i}$ . Therefore,  $V_i(m_i, S'')$  is the same regardless of whether  $i$  has worked or shirked. Agent  $i$ 's expected gross utility conditional on (i) working, (ii) finding no signal, and (iii) producing message  $m_i$ , is

$$V_i^w(\emptyset, m_i) = \sum_{S'' \in \mathcal{S}} V_i(m_i, S'') \Delta_i^\emptyset(S''),$$

whereas his expected gross utility if  $i$  shirks and sends message  $m_i$  is

$$V_i^s(m_i) = \sum_{S'' \in \mathcal{S}} V_i(m_i, S'') \Delta_i(S'')$$

because  $i$  has learned nothing from shirking and thus holds the same belief as his prior belief at the beginning of round  $i$ . Combining these expressions, we get

$$V_i^w(\emptyset, m_i) - V_i^s(m_i) = \sum_{S'' \in \mathcal{S}} V_i(m_i, S'') (\Delta_i^\emptyset(S'') - \Delta_i(S'')). \quad (27)$$

Since  $V_i(m_i, S'') \in [0, R]$  for all  $m_i$  and  $S''$ , combining (27) with (25) and (26) yields

$$V_i^w(\emptyset, m_i) - V_i^s(m_i) \leq \frac{R \Delta_i(\emptyset) \lambda (1 - f_i^0)}{(1 - f_i^0)(1 - \lambda) + f_i^0}.$$

Since  $\lambda < 1$ , the denominator is bounded below by  $1 - \lambda$ . Since  $\Delta_i(\emptyset) = f_i^0$ , the numerator is bounded above by  $R f_i^0$ . This yields

$$V_i^w(\emptyset, m_i) \leq V_i^s(m_i) + f_i^0 \frac{R}{1 - \lambda} \leq V_i^* + f_i^0 \frac{R}{1 - \lambda},$$

which proves the lemma. Intuitively, this results means that if  $f_i^0$  is negligible relative to  $(1 - \lambda)$ , then  $i$ 's expected gross utility after working and finding nothing cannot be much higher than if  $i$  had shirked, because finding nothing in this case merely reveals that  $i$  was unlucky and otherwise conveys little else information.

### C.3 Proof of Lemma 3

For each  $m_i \in M$ , let  $V_i^w(m_i)$  denote  $i$ 's expected gross utility conditional on working and sending message  $m_i$  and  $M_i^-$  denote the set of messages  $m_i$  after which no  $j > i$  ever works, so that  $M = M_i^+ \cup M_i^-$  and  $M_i^+ \cap M_i^- = \emptyset$ . Letting  $\gamma_i(\tilde{M}_i)$  denote the probability that  $i$  sends a message in  $\tilde{M}_i$  conditional on working and on  $m_1^{i-1}$ , we have:

$$V_i^w = \sum_{m_i \in M_i^-} \gamma_i(m_i) V_i^w(m_i) + \sum_{m_i \in M_i^+} \gamma_i(m_i) V_i^w(m_i). \quad (28)$$

For the first term, note that  $i$ 's expected utility conditional on reporting  $m_i$  and on no  $j > i$  ever producing real evidence does not depend on whether  $i$  worked or shirked: either way, the distribution of the reports  $\{m_j\}_{j>i}$  is independent of the set of signals that remain in the case. Letting, as in the previous lemma,  $V_i^s(m_i)$  denote  $i$ 's expected gross utility conditional on shirking and sending message  $m_i$ , we thus have  $V_i^w(m_i) = V_i^s(m_i)$  for all  $m_i \in M_i^-$ . Since  $V_i^* = \max_{m_i \in M} V_i^s(m_i)$ , the first term in (28) is bounded above by  $\gamma_i(M_i^-) V_i^*$ .

For the second term, we have  $\gamma_i(m_i) = d_i(m_i) + g_i(m_i)$  and

$$\gamma_i(m_i) V_i^w(m_i) \leq d_i(m_i) V_i^w(\emptyset, m_i) + g_i(m_i) R,$$

where we used the fact that  $i$ 's expected gross utility conditional on working, finding a signal, and reporting  $m_i$  is bounded by  $R$ .

Combining these observations yields

$$V_i^w \leq \gamma_i(M_i^-) V_i^* + g_i(M_i^+) R + \sum_{m_i \in M_i^+} d_i(m_i) V_i^w(\emptyset, m_i). \quad (29)$$

If  $\lambda < 1$ , Lemma 2 implies that  $V_i^w(\emptyset, m_i) \leq V_i^* + \frac{f_i^0 R}{1-\lambda}$ . Summing over all  $m_i \in M_i^+$ , we get

$$\sum_{m_i \in M_i^+} d_i(m_i) V_i^w(\emptyset, m_i) \leq d_i(M_i^+) V_i^* + d_i(M_i^+) \frac{f_i^0 R}{1-\lambda}. \quad (30)$$

Since  $\gamma_i(M_i^-) + d_i(M_i^+) \leq \gamma_i(M_i^-) + \gamma_i(M_i^+) = 1$ , combining (29) and (30) proves the lemma when  $\lambda < 1$ .

If  $\lambda = 1$ , using in (29) the fact that  $V_i^w(\emptyset, m_i)$  is bounded above by  $R$  directly proves the lemma.

## C.4 Proof of Proposition 2

The right-hand side of (5) is increasing in  $C$  over the range of  $C$  that satisfy the condition  $F_i^k - CF_i^{k+1} > 0$ . Therefore, if (5) is satisfied for any  $C$  such that  $F_i^k - CF_i^{k+1} > 0$ , it is also satisfied for any  $C' \geq C$  such that  $F_i^k - C'F_i^{k+1} > 0$ . The proposition thus follows if we show the inequality for  $C(\lambda)$ .

First, we show that the claim holds if  $\beta_i = 0$ . In this case,  $i$  must shirk with probability 1: if not, Lemma 1 implies that  $S_i$  is nonempty with positive probability and, hence, that  $\beta_i > 0$ . Since  $i$  shirks with probability 1, there is no belief update between rounds  $i$  and  $i + 1$ . Therefore,  $F_i^k = F_{i+1}^k(m_i)$  for any message  $m_i$  that  $i$  sends in equilibrium. The right-hand side of (5) is thus equal to zero and (5) is satisfied.

Now suppose that  $\beta_i > 0$  or, equivalently, that  $\gamma_i > 0$ :  $i$  works with positive probability. We consider two cases, distinguished by the magnitude of the probability  $g_i(M_i^+)$  that  $i$  finds a signal and sends a message in  $M_i^+$  conditional on working and on history  $m_1^{i-1}$ .

**Case 1:**  $g_i(M_i^+) \geq \frac{c}{2R}$ . By definition, we have

$$\beta_i = \gamma_i \lambda (1 - f_i^0),$$

which implies that  $\beta_i \leq \gamma_i$ , and

$$\beta_i(M_i^+) = \gamma_i g_i(M_i^+)$$

where  $\beta_i(\tilde{M}_i)$  is the probability that  $i$  discovers a signal and sends a message in  $\tilde{M}_i$ . Since  $g_i(M_i^+) \geq c/2R$ , this implies that

$$\beta_i \leq \gamma_i \leq \frac{2R}{c} \beta_i(M_i^+). \quad (31)$$

Therefore, the desired inequality (5) will follow for  $C = 2R/c$  if we prove that  $\beta_i(M_i^+)$  is bounded above by  $\frac{\mathbb{E}_i \left[ \frac{(F_i^k - F_{i+1}^k(m_i)) \mathbf{1}_{m_i \in M_i^+}}{F_i^k - CF_i^{k+1}} \right]}$ .

For each  $m_i$  that  $i$  may send in equilibrium and  $k \geq 1$ , Bayesian updating implies that

$$F_{i+1}^k(m_i) = \frac{F_i^k(\alpha_i(m_i) + \gamma_i(1 - \lambda)\delta_i(m_i)) + \Phi(m_i)}{\alpha_i(m_i) + \gamma_i((1 - F_i^1) + F_i^1(1 - \lambda))\delta_i(m_i) + \beta_i(m_i)} \quad (32)$$

where the following probabilities are defined conditional on  $m_1^{i-1}$ :

- $\alpha_i(\tilde{M}_i)$ : probability that  $i$  shirks and sends a message in  $\tilde{M}_i$ ;

- $\delta_i(\tilde{M}_i)$ : probability that  $i$  sends a message in  $\tilde{M}_i$  conditional on working *and* finding no signal;<sup>59</sup>
- $\Phi(m_i)$  is the probability that (i)  $i$  works, (ii)  $i$  discovers a signal, (iii)  $i$  sends report  $m_i$ , and (iv) there remain at least  $k$  signals at the beginning of round  $i + 1$ .

Let  $p_i(m_i)$  denote the probability that  $i$  produces report  $m_i$  conditional on  $m_1^{i-1}$ :  $p_i(m_i)$  is the denominator of (32). Rearranging (32) and simplifying, we have

$$F_i^k(\beta_i(m_i) + \gamma_i(1 - F_i^1)\lambda\delta_i(m_i)) = (F_i^k - F_i^{k+1}(m_i))p_i(m_i) + \Phi(m_i) \quad (33)$$

Since  $\gamma_i(1 - F_i^1)\lambda\delta_i(m_i) \geq 0$ , summing the previous equation over  $m_i \in M_i^+$  yields

$$F_i^k\beta_i(M_i^+) \leq \mathbb{E}_i \left[ (F_i^k - F_{i+1}^k(m_i)) \mathbb{1}_{m_i \in M_i^+} \right] + \sum_{m_i \in M_i^+} \Phi(m_i). \quad (34)$$

Since  $\Phi(m_i) = \mathbb{E}_i \left[ \mathbb{1}_{|S_i| \geq k+1} \mathbb{1}_{i \text{ works, discovers a signal, and reports } m_i} \right]$ , we have

$$\begin{aligned} \sum_{m_i \in M_i^+} \Phi(m_i) &\leq \sum_{m_i \in M_i^+} \mathbb{E}_i \left[ \mathbb{1}_{|S_i| \geq k+1} \mathbb{1}_{i \text{ works and reports } m_i} \right] \\ &= \mathbb{E}_i \left[ \mathbb{1}_{|S_i| \geq k+1} \mathbb{1}_{i \text{ works and reports } m_i \in M_i^+} \right] \\ &\leq \mathbb{E}_i \left[ \mathbb{1}_{|S_i| \geq k+1} \mathbb{1}_{i \text{ works}} \right] \\ &= F_i^{k+1}\gamma_i, \end{aligned} \quad (35)$$

noting, for the last equality, that the event that  $i$  works, which has probability  $\gamma_i$ , depends only on  $m_1^{i-1}$  and is thus independent of the event  $\{|S_i| \geq k + 1\}$  conditional on  $m_1^{i-1}$ .

Combining this with (34) yields

$$F_i^k\beta(M_i^+) \leq \mathbb{E}_i \left[ (F_i^k - F_{i+1}^k(m_i)) \mathbb{1}_{m_i \in M_i^+} \right] + F_i^{k+1}\gamma_i. \quad (36)$$

Since  $g_i(M_i^+) \geq c/2R$ , we have  $\beta(M_i^+) = \gamma_i g_i(M_i^+) \geq \gamma_i c/2R$ . Inequality (36) then yields

$$\beta(M_i^+) \leq \frac{1}{F_i^k - 2R/cF_i^{k+1}} \mathbb{E}_i \left[ (F_i^k - F_{i+1}^k(m_i)) \mathbb{1}_{m_i \in M_i^+} \right]. \quad (37)$$

Combining this with (31) yields (5) for  $C(1) = 2R/c$ .<sup>60</sup> Since  $C(\lambda) \geq C(1)$  for all  $\lambda \in (0, 1]$ , the monotonicity noted at the beginning of the proof yields the desired conclusion for  $C(\lambda)$ .

<sup>59</sup>Note that  $\delta_i(\tilde{M}_i) \geq d_i(\tilde{M}_i)$ , where  $d_i(\tilde{M}_i)$  was defined before Lemma 3.

<sup>60</sup>Note that the proposition's assumption that  $F_i^k - C(\lambda)F_i^{k+1} > 0$  implies that  $F_i^k - \frac{2R}{c}F_i^{k+1} > 0$  since  $C(\lambda) \geq C(1)$  regardless of  $\lambda$ .

**Case 2:**  $g_i(M_i^+) < \frac{c}{2R}$ . We prove that  $\gamma_i$  is bounded above by the right-hand side of (5) for  $C = C(\lambda)$ . Since  $\gamma_i \geq \beta_i$ , this will yield the desired conclusion.

Intuitively, in Case 2 the probability of discovering a signal that, together with  $i$ 's equilibrium message strategy, triggers subsequent work is too low to incentivize  $i$  to work. The only way of incentivizing  $i$  to work is therefore for him to signal by his message that he found nothing through work. For this to happen, the probability  $f_i^0$  that there remains no evidence must be high enough. We will use this fact to obtain a bound on  $\gamma_i$ .

From Lemma 3, if  $g_i(M_i^+) < c/2R$ ,  $i$ 's utility from working is bounded above by

$$U_i^w = V_i^w - c \leq V_i^* + d_i(M_i^+) \frac{f_i^0 R}{1 - \lambda} - \frac{c}{2}$$

if  $\lambda < 1$ , and by

$$U_i^w \leq V_i^* + d_i(M_i^+) R - \frac{c}{2}$$

if  $\lambda = 1$ . Therefore, working is optimal only if  $d_i(M_i^+) f_i^0 \geq c(1 - \lambda)/2R$  when  $\lambda < 1$  and only if  $d_i(M_i^+) \geq c/2R$  when  $\lambda = 1$ .

Summing (33) over  $M_i^+$  and using (35) and  $f_i^0 = 1 - F_i^1$  yields

$$F_i^k \beta_i(M_i^+) + \gamma_i F_i^k \delta_i(M_i^+) f_i^0 \lambda \leq \mathbb{E}_i \left[ (F_i^k - F_{i+1}^k(m_i)) 1_{m_i \in M_i^+} \right] + F_i^{k+1} \gamma_i. \quad (38)$$

For  $\lambda < 1$ , we have  $d_i(M_i^+) f_i^0 \geq c(1 - \lambda)/2R$ . Since  $\delta_i(m_i) \geq d_i(m_i)$  for all  $m_i$  (by definition of these variables) and  $F_i^k \beta_i(M_i^+) \geq 0$ , (38) implies that

$$\gamma_i \leq \frac{\mathbb{E}_i \left[ (F_i^k - F_{i+1}^k(m_i)) 1_{m_i \in M_i^+} \right]}{F_i^k c \lambda (1 - \lambda) / 2R - F_i^{k+1}}.$$

Multiplying the numerator and denominator by  $C(\lambda)$  yields the result.

For  $\lambda = 1$ , we have  $d_i(m_i) = \delta_i(m_i) f_i^0$  for all  $m_i$  and, hence,  $\delta_i(M_i^+) f_i^0 \lambda = d_i(M_i^+)$ , which is greater than  $c/2R$  as noted earlier. Therefore, (38) implies that

$$\gamma_i \leq \frac{\mathbb{E}_i \left[ (F_i^k - F_{i+1}^k(m_i)) 1_{m_i \in M_i^+} \right]}{F_i^k (c/2R) - F_i^{k+1}}.$$

Multiplying the numerator and the denominator by  $C(1)$  yields the result. ■

## C.5 Proof of Lemma 5

(i) Let  $S'_i = \{s_i : F_i^{k,w}(s_i) > F_i^k\}$ , and, for each  $s_i$ , let  $\gamma'_i(s_i)$  (resp.  $\beta'_i(s_i)$ ) denote the probability that  $i$  discovers  $s_i$  given that he works (resp., the probability that  $i$  works and discovers  $s_i$ ). Also let  $\gamma'_i(s_i|k+1)$  (resp.  $\beta'_i(s_i|k+1)$ ) denote the same probabilities conditional on  $|S_i| \geq k+1$ .

We have for all  $s_i$

$$F_i^{k,w}(s_i) = \frac{F_i^{k+1}\gamma'_i(s_i|k+1)}{\gamma'_i(s_i)} = \frac{F_i^{k+1}\beta'_i(s_i|k+1)}{\beta'_i(s_i)}$$

where the second equality comes from  $\beta'_i(s_i) = \gamma_i\gamma'_i(s_i)$  and  $\beta'_i(s_i|k+1) = \gamma_i\gamma'_i(s_i|k+1)$ . Therefore,  $F_i^{k,w}(s_i) > F_i^k$  only if  $\beta'_i(s_i) < \beta'_i(s_i|k+1)F_i^{k+1}/F_i^k$ .

We have  $\sum_{s_i \in S'_i} F_i^{k+1}\beta'_i(s_i|k+1) \leq \gamma_i F_i^{k+1}$ . Therefore, the probability  $\beta'_i(S'_i)$  that  $i$  works and finds a signal in  $S'_i$  satisfies

$$\beta'_i(S'_i) = \sum_{s_i \in S'_i} \beta'_i(s_i) \leq \gamma_i \frac{F_i^{k+1}}{F_i^k}.$$

Since  $\beta'_i(S'_i) = \gamma_i\gamma'_i(S'_i)$ , we have

$$\gamma'_i(S'_i) \leq \frac{F_i^{k+1}}{F_i^k} \leq \frac{\varepsilon}{\hat{F}^k},$$

since  $F_i^{k+1} \leq \varepsilon$  and  $F_i^k \geq \hat{F}^k$ . From (6), the right-hand side is bounded above by  $\frac{\hat{F}^k}{2G^2}$ .

For any  $m_i$ , let  $q(m_i)$  denote the probability, conditional on  $i$  working and sending report  $m_i$ , that  $i$  has discovered a signal  $s_i \in S'_i$ , and let  $\sigma(s_i|m_i)$  denote the probability that  $i$  discovered  $s_i$  given that he worked and reported  $m_i$ . We also let  $s_i = \emptyset$  denote the event that  $i$  did not find anything,  $\sigma(\emptyset|m_i)$  denote the probability that  $i$  found nothing given that he worked and reported  $m_i$ ,  $F_i^{k,w}(\emptyset)$  denote the probability that there at least  $k$  signals conditional on  $i$  working and finding nothing. We have

$$\begin{aligned} F_i^{k,r}(m_i) &= \sum_{s_i \in S_i} \sigma(s_i|m_i) F_i^{k,w}(s_i) \\ &= \sum_{s_i \in S'_i} \sigma(s_i|m_i) F_i^{k,w}(s_i) + \sigma(\emptyset|m_i) F_i^{k,w}(\emptyset) + \sum_{s_i \neq \emptyset, s_i \in S_i \setminus S'_i} \sigma(s_i|m_i) F_i^{k,w}(s_i) \end{aligned}$$

By construction,  $F_i^{k,w}(s_i) \leq F_i^k$  for all  $s_i$  in the last term. Moreover, we have  $F_i^{k,w}(\emptyset) \leq F_i^k$ , as is easily checked.<sup>61</sup> intuitively, finding nothing always increases the probability that there

<sup>61</sup>Formally, for  $k \geq 1$ , we have  $F_i^{k,w}(\emptyset) = \frac{F_i^k(1-\lambda)}{(1-F_i^1)+F_i^1(1-\lambda)} = F_i^k \frac{1-\lambda}{(1-F_i^1)+F_i^1(1-\lambda)} \leq F_i^k$ .

are no signals remaining to be found. Finally, the first term is bounded above by  $\sigma(S'_i|m_i) = q(m_i)$ . Therefore,

$$\begin{aligned} F_i^{k,r}(m_i) \geq (1 + \eta)F_i^k &\Rightarrow q(m_i) + (1 - q(m_i))F_i^k \geq (1 + \eta)F_i^k \\ &\Rightarrow q(m_i) \geq \eta F_i^k. \end{aligned}$$

To conclude, note that

$$\sum_{m_i} \gamma_i(m_i)q(m_i) = \Pr(s_i \in S'_i | m_1^{i-1}, i \text{ works}) \leq \frac{\hat{F}^k}{2G^2}$$

The left-hand side is bounded below by  $\gamma(N_i)\eta F_i^k$ . Since  $F_i^k \geq \hat{F}^k$ , this implies that

$$\gamma(N_i) \leq \frac{1}{2\eta G^2}.$$

(ii)  $F_{i+1}^k(m_i)$  is a convex combination<sup>62</sup> of  $F_i^k$  and  $F_i^{k,r}(m_i)$ . This implies that  $F_{i+1}^k(m_i) \leq (1 + \eta)F_i^k$  for all  $m_i \notin N_i$ . From (14), this further implies that

$$\pi_{i+1}(m_i) \leq \frac{2Cg}{\hat{F}^k} ((1 + \eta)F_i^k - \mathbb{E}_{i+1}[F_i^k])$$

for all  $m_i \notin N_i$ .

## C.6 Proof of Lemma 6

(i) If  $m_i \notin T_i$ , we have  $F_{i+1}^{k+1} \leq \sqrt{G}\varepsilon$ . Using this inequality in Lemma 4 instead of  $F_i^{k+1} \leq \varepsilon$  and repeating the argument of that lemma applied to round  $i + 1$ , we conclude that  $i + 1$  assigns probability at least  $1 - 1/\sqrt{G}$  to  $\mathcal{A}$  whenever  $m_i \notin T_i$ .

(ii) Let  $F_i^{k+1,r}(m_i)$  denote the probability that there are at least  $k + 1$  signals left at the beginning of round  $i + 1$  conditional on  $i$  working and reporting  $m_i$ .  $F_{i+1}^{k+1}$  is a convex combination of  $F_i^{k+1}$  and  $F_i^{k+1,r}(m_i)$ . This, together with the fact that  $F_i^{k+1} \leq \varepsilon$  and the

---

<sup>62</sup> $F_{i+1}^k(m_i)$  is the probability that  $i + 1$  assigns to there being at least  $k$  signals left upon observing  $m_i$ . If  $i + 1$  knew that  $i$  didn't work and simply sent message  $m_i$ , this belief should be  $F_i^k$  since  $m_i$  conveys no additional information. And if  $i + 1$  knew that  $i$  produced  $m_i$  through working and then reporting  $m_i$ , his updated belief should be  $F_i^{k,r}(m_i)$ . Since  $i + 1$  doesn't observe  $i$ 's action, in general  $F_{i+1}^k(m_i)$  is a convex combination of these two posteriors, where the weights corresponds to the probability assigned by  $i + 1$  to  $i$  fabricating or working conditional on observing  $m_i$ . This fact is straightforward to check using Bayesian updating.



definition of  $T_i$ , shows that  $m_i \in T_i$  only if  $F_i^{k+1,r}(m_i) \geq \sqrt{G}\varepsilon$ . Let  $T'_i$  denote the set of messages  $m_i$  for which the last inequality holds. As noted,  $T_i \subset T'_i$ .

Since  $i$ 's prior probability that there are at least  $k+1$  signals left is  $F_i^{k+1} \leq \varepsilon$ , the law of iterated expectations implies that the probability  $\bar{F}_i^{k+1,r}(m_i)$  that there were at least  $k+1$  signals left at the beginning of round  $i$  conditional on  $i$  working and finding  $m_i$  must satisfy

$$\mathbb{E}_i[\bar{F}_i^{k+1,r} | \text{working}] = \sum_{m_i \in M_i} \gamma_i(m_i) \bar{F}_i^{k+1,r}(m_i) = F_i^{k+1} \leq \varepsilon.$$

Using Markov's inequality, this implies that  $\Pr(m_i : \bar{F}_i^{k+1,r}(m_i) \geq \sqrt{G}\varepsilon | i \text{ works}) \leq \frac{\varepsilon}{\sqrt{G}\varepsilon} = 1/\sqrt{G}$ . Since also  $\bar{F}_i^{k+1,r}(m_i) \geq F_i^{k+1,r}(m_i)$ , we get  $\gamma_i(T'_i) \leq 1/\sqrt{G}$ . Since  $T_i \subset T'_i$ , this shows that  $\gamma_i(T_i) \leq 1/\sqrt{G}$ .

## D Proof of Theorem 3

Without loss of generality, we assume once more that agents' gross utility functions  $\{V_i\}_{i \in \mathbb{N}}$  all take values in  $[0, R]$ .

Conditional on the history up to round  $i$ , the number of signals discovered until round  $i$  is a random variable whose support depends on  $m_1^{i-1}$  and the set of past witnesses. Let  $q_i$  denote the lower bound of this support. Put differently,  $q_i$  is the smallest number of signals that have been surely discovered by round  $i$ , and it is equal to the number of past witnesses plus the number of past investigators whose equilibrium strategies and messages imply that they have surely discovered signals given the history up to round  $i$ . Let  $\mathcal{Q}_i$  denote the event that the number of signals discovered up to round  $i$  is exactly equal to  $q_i$ , and  $\hat{F}_i^k$  denote the probability that  $|S_i| \geq k$  conditional on  $\mathcal{Q}_i$  and  $q_i$ .

We observe that (i)  $q_i$  is nondecreasing along any equilibrium path and is strictly increasing whenever a witness arrives, (ii) given the construction of  $S$ , the distribution of  $S_i$  conditional on  $m_1^{i-1}$  and  $\mathcal{Q}_i$  is only a function of  $q_i$ , and (iii)  $\hat{F}_1^k = F_1^k$  for all  $k$ .

We will prove that there exist strictly positive thresholds  $\{F^k\}_{k \geq 1}$  such that an informative continuation equilibrium exists in round  $i$  only if  $\hat{F}_i^k \geq F^k$  for all  $k \geq 1$ . Applied to  $i = 1$ , this result implies Theorem 3. The proof uses the following lemma.

**Lemma 7** *Under Assumption 1, the following inequalities hold: (i)  $F_i^k \leq \hat{F}_i^k$  for all  $i, k \geq 1$  and (ii) Path by path,  $\hat{F}_j^k \leq \hat{F}_i^k$  for all  $j \geq i$  and  $k \geq 1$ .*

*Proof.* Part (i): Let  $r_i$  denote the number of signals discovered before round  $i$ . We have  $r_i \geq q_i$  and

$$\begin{aligned}
F_i^k &= \sum_{r \in \{q_i, \dots, i-1\}} \Pr(r_i = r \mid m_1^{i-1}) \Pr(\tilde{K} \geq r + k \mid \tilde{K} \geq r) \\
&\leq \sum_{r \in \{q_i, \dots, i-1\}} \Pr(r_i = r \mid m_1^{i-1}) \Pr(\tilde{K} \geq q_i + k \mid \tilde{K} \geq q_i) \\
&= \sum_{r \in \{q_i, \dots, i-1\}} \Pr(r_i = r \mid m_1^{i-1}) \hat{F}_i^k \\
&= \hat{F}_i^k.
\end{aligned}$$

The first equality comes from the independence of  $\tilde{K}$  from  $S^\infty$ :  $r_i$  is a sufficient statistic for  $\tilde{K}$  given all the information produced before round  $i$ , and only to the extent that it reveals that  $\tilde{K} \geq r_i$ . The inequality comes from the increasing hazard rate condition, which implies that for any  $k \geq 0$ ,  $\Pr(\tilde{K} \geq k + q \mid \tilde{K} \geq q)$  is non-increasing in  $q$ .<sup>63</sup> The second equality is due to the equality  $\hat{F}_i^k = \Pr(\tilde{K} \geq k + q_i \mid \tilde{K} \geq q_i)$ , which again comes from the independence of  $\tilde{K}$  and  $S^\infty$ : the only relevant information about  $\tilde{K}$  conditional on  $\mathcal{Q}_i$  is the number of signals  $q_i$  discovered by round  $i$ .

Part (ii) For any  $j \geq i$ , we have  $q_j \geq q_i$ . As explained in the proof of Part (i), we have for all  $k \geq 1$

$$\begin{aligned}
\hat{F}_j^k &= \Pr(\tilde{K} \geq k + q_j \mid \tilde{K} \geq q_j) \\
&\leq \Pr(\tilde{K} \geq k + q_i \mid \tilde{K} \geq q_i) \\
&= \hat{F}_i^k,
\end{aligned}$$

where the inequality comes from the increasing hazard rate property (see Footnote 63). ■

The existence of thresholds  $\{\underline{F}^k\}_{k \geq 1}$  in Theorem 3 is proved by induction on  $k$ . We start with the base case  $k = 1$  and then prove the induction step.

---

<sup>63</sup>See, e.g., Barlow et al. (1963, p. 379). In brief, a random variable  $X$  with distribution  $F$  has the increasing hazard rate property if and only if the survival distribution  $\bar{F} = 1 - F$  is log-concave over its domain (which could be discrete or continuous). This property implies, as is easily checked, that for any  $p, q \geq 0$ ,  $\Pr(X \geq p + q \mid X \geq p) = \bar{F}(p + q)/\bar{F}(p)$  is nonincreasing in  $p$ .

## D.1 Proof for $k = 1$

Suppose that  $\hat{F}_i^1 < c/2R$ . Combining the two parts of Lemma 7, this implies that  $F_j^1 < c/2R$  for all  $j \geq i$ . Lemma 1 still applies: no investigator  $j \geq i$  works because the probability that he finds something is too small to justify the cost of effort, given the maximal reward  $R$ .

We now show that if  $\hat{F}_i^1$  lies below another threshold, smaller than  $c/2R$ , witnesses provide no informative message, either.

We will use the following lemma:

**Lemma 8** *Consider  $L \geq 2$  pairwise independent random variables  $\{Y_\ell\}$  with non-atomic distributions over  $\mathbb{R}$  and densities  $f_\ell$  that are bounded above by  $\bar{f}$ . For any  $\varepsilon \geq 0$ , let  $E_\varepsilon^L$  denote the event that there exist  $\ell, \ell' \leq L$  such that  $|Y_\ell - Y_{\ell'}| \leq \varepsilon$ . Then*

$$\Pr(E_\varepsilon^L) \leq L(L-1)\bar{f}\varepsilon.$$

*Proof.* The result is proved by induction on  $L \geq 2$ . For  $L = 2$ , we have

$$\Pr(|Y - Y'| \leq \varepsilon) = \int_{\mathbb{R}} f_Y(x) F_{Y'}[x - \varepsilon, x + \varepsilon] dx \leq 2\varepsilon \bar{f} \int_{\mathbb{R}} f_Y(x) dx = 2\varepsilon \bar{f}.$$

Now suppose that the claim holds for  $L-1$ . Notice that the event  $E_\varepsilon^L$  is the union of  $L$  events: the event  $E_\varepsilon^{L-1}$  involving the first  $L-1$  random variables, and, for each  $\ell \leq L-1$ , the event  $E^{\ell,L}$  that the  $L^{\text{th}}$  random variable lies within  $\varepsilon$  of the  $\ell^{\text{th}}$  random variable. Therefore,

$$\begin{aligned} \Pr(E_\varepsilon^L) &\leq \Pr(E_\varepsilon^{L-1}) + \sum_{\ell \leq L-1} \Pr(|Y_\ell - Y_L| \leq \varepsilon) \\ &\leq (L-1)(L-2)\bar{f}\varepsilon + (L-1) \times 2\varepsilon \bar{f} \\ &= L(L-1)\bar{f}\varepsilon, \end{aligned}$$

where the second inequality comes from the induction hypothesis and the fact that  $\Pr(|Y_\ell - Y_L| \leq \varepsilon) \leq 2\bar{f}\varepsilon$ , as shown in the first step of the induction for  $L = 2$ .  $\blacksquare$

If  $i$  is a witness, we will use the following notation:

- $\beta_i$ : probability that  $i$  produces an informative message given  $m_1^{i-1}$ ;
- $M_i^+$ : set of messages  $m_i$  that are followed by an informative continuation equilibrium;
- $\Pr_i(M_i^+)$ : probability that  $i$  produces a message in  $M_i^+$  given  $m_1^{i-1}$ ;

- $\gamma_i(m_i)$ : probability that  $i$  sends  $m_i$  given  $m_1^{i-1}$ ;
- $\gamma_i(m_i|s_i)$ : probability that  $i$  sends  $m_i$  after observing  $s_i$ .

**Lemma 9** *There is a threshold  $\underline{F}^1 \in (0, c/2R)$  such if  $i$  a witness, then  $\beta_i > 0$  only if  $\hat{F}_i^1 \geq \underline{F}^1$ .*

*Proof.* Recall that if  $i$  is a witness, his message  $m_i$  is informative (in equilibrium) if it is statistically dependent of  $S$  conditional on  $m_1^{i-1}$ . Given a vector  $\epsilon = (\epsilon(m_i) : m_i \in M)$  of shocks, say that  $i$ 's message is  $\epsilon$ -informative if whenever  $i$  has the preference shock  $\epsilon$  as defined by (4), there exist two signals  $s_i \neq s'_i$  such that the equilibrium distributions of  $m_i$  conditional on  $i$  getting signals  $s_i$  versus  $s'_i$  are different across these two signals.

The following observation is straightforward to prove.

**OBSERVATION 1**  *$i$ 's message is informative if and only if the set of preference shocks  $\epsilon$  for which  $i$ 's message is  $\epsilon$ -informative has positive probability.*

For any equilibrium and history up to some round  $i$  in which  $i$  is a witness, let  $\nu_i$  denote the probability that  $i$ 's preference shock  $\epsilon_i$  is such that  $i$ 's message is  $\epsilon_i$ -informative.

For any  $F < c/2R$ , let  $\nu(F)$  denote the supremum of  $\nu_i$  over all witness rounds  $i$  of all equilibria such that  $\hat{F}_i^1 \leq F$ . We will show that  $\nu(F) = 0$  for all  $F$  below some strictly positive threshold.

Consider such an equilibrium. For any witness round  $i$  and message  $m_i$ , let  $z(m_i)$  denote the probability that at least some  $j > i$  produces an informative message following message  $m_i$ .

$i$ 's expected utility if he receives signal  $s_i$  and sends message  $m_i$  is given by

$$U_i(m_i; s_i) = z(m_i)\mathbb{E}_i[V_i(m) | s_i, m_1^i] + (1 - z(m_i))V_i(m_i) + \epsilon_i(m_i) \quad (39)$$

where

$$V_i(m_i) = \mathbb{E}_i[V_i(m) | m_1^i, \text{no } j > i \text{ produces an informative message}].$$

Notice that  $V_i(m_i)$  does not depend on the content of signal  $s_i$  since this signal is payoff irrelevant whenever no  $j > i$  produces an informative message.

Given  $\epsilon_i$ ,  $i$  sends an informative signal only if there exist  $m_i \neq m'_i$  and signals  $s_i \neq s'_i$  such that

$$U_i(m_i; s_i) \geq U_i(m'_i; s_i)$$

and

$$U_i(m'_i; s'_i) \geq U_i(m_i; s'_i)$$

From (39) and the fact that  $V_i(m) \in [0, R]$ , this is possible only if:

$$|\epsilon_i(m_i) + V_i(m_i) - \epsilon_i(m'_i) - V_i(m'_i)| \leq R(z(m_i) + z(m'_i)).$$

The random variables  $Y_\ell = \epsilon_i(m_\ell) + V_i(m_\ell)$  satisfy the assumptions of Lemma 8. Letting  $|M|$  denote the cardinality of the message space  $M$ , we thus have

$$\Pr(i \text{ sends an informative message}) \leq |M|^2 \bar{f} R \sup_{m_i, m'_i \in M} (z(m_i) + z(m'_i)) \quad (40)$$

Since no investigator  $j \geq i$  works when  $\hat{F}_i^1 < c/2R$ , we have for any message  $m_i$ :

$$z(m_i) \leq \sum_{j \geq 1} \Pr(\text{there are } j \text{ witnesses in the sequence after round } i) j \nu(F). \quad (41)$$

Indeed, by definition of  $\nu(F)$  and the fact that  $\hat{F}_j^1 \leq F$  for all  $j \geq i$ , a witness provides an informative signal with probability at most  $\nu(F)$ . Thus  $j\nu(F)$  is an upper bound on the probability that at least one witness provides an informative signal given there are  $j$  such witnesses.

The probability that at least  $j$  witnesses come after round  $i$  is bounded above by  $F^j$ , where  $j$  is an exponent (not a superscript): To show this for  $j = 1$ , notice that by Lemma 7,  $\hat{F}_{i+1}^1 \leq F$ , and the probability that there is at least one witness after round  $i$  is bounded above by the probability  $F_{i+1}^1$  that there is at least one more signal, which is less than  $\hat{F}_{i+1}^1$  again by Lemma 7. For  $j = 2$ , note that conditional on the first witness arriving, Lemma 7 implies that the probability that a second witness arrives is again bounded by  $F$  since the probability that there remains another signal is bounded by  $F$ , and a witness can arise only if such a signal exists. By induction, this shows that the probability of having at least  $j$  witnesses and, hence, the probability of having exactly  $j$  witnesses, are bounded above by  $F^j$ . Combining this with (41) and using the standard formula  $\sum_{j \geq 1} jx^j = x/(1-x)^2$  for all  $x \in (0, 1)$ , we get for all  $m_i$ :

$$z(m_i) \leq \sum_{j \geq 1} F^j j \nu(F) = \frac{F \nu(F)}{(1-F)^2}.$$

Combining this with (40), we obtain

$$\Pr(i \text{ sends an informative message} \mid \hat{F}_i^1 \leq F) \leq 2|M|^2 \bar{f} R \nu(F) \frac{F}{(1-F)^2}. \quad (42)$$

Taking the supremum of the left-hand side over all witness rounds  $i$  and equilibria such that  $\hat{F}_i^1 \leq F$ , we obtain

$$\nu(F) \leq |M|^2 \bar{f} R \frac{2F}{(1-F)^2} \nu(F).$$

For  $2F/(1-F)^2 < 1/|M|^2 \bar{f} R$ , this relation is possible only if  $\nu(F) = 0$ . Since the function  $F \mapsto F/(1-F)^2$  is increasing on  $[0, 1)$  and starts at zero, we conclude that there exists a threshold  $\underline{F}^1 > 0$  such that  $\nu(F) = 0$  for all  $F \leq \underline{F}^1$ .  $\blacksquare$

## D.2 Induction Step

Suppose that there exist strictly positive thresholds  $\{\underline{F}^{k'}\}_{k' \in \{1, \dots, k\}}$  such that a continuation equilibrium starting in round  $i$  is informative only if  $\hat{F}_i^{k'} \geq \underline{F}^{k'}$  for all  $k' \leq k$ . We will show that a similar condition holds for  $k+1$ . The proof works by contradiction: we will suppose that for all  $\varepsilon \in (0, 1)$ , there exists an informative continuation equilibrium such that  $\hat{F}_i^{k+1} \leq \varepsilon$  and obtain an impossibility for  $\varepsilon$  small enough.

Consider any  $\varepsilon < \underline{F}^k \times \underline{F}^1$  and any informative continuation equilibrium starting in round  $i$  such that  $\hat{F}_i^{k+1} \leq \varepsilon$ . Then, all continuation equilibria become uninformative as soon as some witness  $j \geq i$  arrives. To see this, suppose that some witness arrives in round  $j \geq i$ . We have for any message  $m_j$  sent by this witness:

$$\begin{aligned} \hat{F}_{j+1}^k(m_j) &= \Pr(\tilde{K} \geq k + q_j + 1 \mid \tilde{K} \geq q_j + 1) \\ &= \frac{\Pr(\tilde{K} \geq k + q_j + 1)}{\Pr(\tilde{K} \geq q_j + 1)} \\ &\leq \frac{\Pr(\tilde{K} \geq k + q_i + 1)}{\Pr(\tilde{K} \geq q_i + 1)} \\ &= \frac{\Pr(\tilde{K} \geq k + q_i + 1)}{\Pr(\tilde{K} \geq q_i)} \times \frac{\Pr(\tilde{K} \geq q_i)}{\Pr(\tilde{K} \geq q_i + 1)} \\ &= \frac{\hat{F}_i^{k+1}}{\hat{F}_i^1} \end{aligned}$$

where the inequality comes for the monotone hazard rate assumption (see Footnote 63) and the fact that  $q_j \geq q_i$ . Since the continuation equilibrium from round  $i$  is informative, we must have  $\hat{F}_i^1 \geq \underline{F}^1$ . Therefore,  $\hat{F}_i^{k+1} \leq \varepsilon < \underline{F}^k \underline{F}^1$ , which implies that

$$\hat{F}_{j+1}^k < \underline{F}^k$$

and shows by induction hypothesis that all continuation equilibria are uninformative from round  $j + 1$  onwards.

Moreover, the first witness,  $j$ , knowing that continuation equilibria are uninformative regardless of his message, has no incentive to send an informative message.<sup>64</sup>

Therefore, if  $\hat{F}_i^{k+1} \leq \varepsilon$ , the only agents who may send informative messages are the investigators arriving between round  $i$  and the arrival of the first witness. The situation is therefore almost identical to the setting of Theorem 2, in the absence of witnesses, except that any sequential learning activity is interrupted at the apparition of the first witness. We know from Theorem 2 that such equilibria can be informative only if  $F_i^{k+1}$  exceeds the  $(k + 1)^{th}$ -threshold given by Theorem 2, which we denote here  $\tilde{F}^{k+1}$ . Since  $F_i^{k+1} \leq \hat{F}_i^{k+1}$  by Part (i) of Lemma 7, we conclude, letting  $\underline{F}^{k+1} = \min\{\tilde{F}_i^{k+1}, \underline{F}^k \underline{F}^1\}$ , that no informative continuation equilibrium exists in round  $i$  if  $\hat{F}_i^{k+1} \leq \underline{F}^{k+1}$ . This concludes the induction step. ■

## E Proof of Theorem 4

The proof is similar to the proof of Theorem 3. In particular, the preliminaries of Section D and Lemma 7 still apply. Lemma 9 starts identically. Instead of  $z(m_i)$ , we use  $z_i^w$ , which denotes the probability that there exists some agent  $j > i$  who works conditional on the history up to round  $i$  and on the fact that analyst  $i$  works. We redefine  $\nu(F)$  to denote the supremum, over all analyst rounds such that  $\hat{F}_i^1 \leq F$ , of the probability that  $i$  works. As in Lemma 9, we have

$$z_i^w \leq F\nu(F)/(1 - F)^2. \quad (43)$$

Consider  $i$ 's incentive to work: if  $i$  works, his expected utility equals

$$U_i^w = z_i^w V_i^w + (1 - z_i^w) \underline{V}_i - c_i \quad (44)$$

where  $V_i^w$  is  $i$ 's expected gross utility conditional on  $i$  working and some  $j > i$  also working, and  $\underline{V}_i$  is  $i$ 's expected gross utility conditional on  $i$  being the last person to work. If  $i$  shirks and chooses the message  $m_i$  optimally, his expected utility is

$$U_i^s = \sum_{m_i \in M} \mathbb{E}_i[V_i(m)|m_1^i]. \quad (45)$$

---

<sup>64</sup>Formally, the argument is similar to the proof of Lemma 9 except that here  $z(m_i) = 0$  regardless of the message.

Since  $V_i^w \leq R$ ,  $U_i^s \geq V_i$  and (recalling the normalization, used in previous proofs, that  $V_i(m) \geq 0$  for all  $m$ )  $U_i^s \geq 0$ , we deduce from (44) and (45) that

$$U_i^w \leq U_i^s + z_i^w R - c_i.$$

Therefore,  $i$  works only if  $c_i \leq z_i^w R$  and the probability that  $i$  works at the beginning of round  $i$  is bounded above by  $H(z_i^w R)$ . Letting  $\bar{h}$  denote an upper bound on the density of  $H$  near zero and using (43), we conclude that

$$\Pr(i \text{ works}) \leq \bar{h} R F \nu(F) / (1 - F)^2.$$

Taking the supremum over all analysts  $i$  such that  $\hat{F}_i^1 \leq F$ , we obtain

$$\nu(F) \leq \bar{h} R F \nu(F) / (1 - F)^2.$$

The rest of the proof is identical to Sections D.1 and D.2.