

# Economics Department University of Southampton Southampton SO17 1BJ, UK

Discussion Papers in Economics and Econometrics

**Optimal Culpability in Research Teams** 

Erika Domotor (University of Cyprus)

Zacharias Maniadis

(University of Cyprus, University of Southampton)

Nikolas Tsakas (University of Cyprus)

No. 2406

This paper is available on our website http://www.southampton.ac.uk/socsci/economics/research/papers

ISSN 0966-4246

# **Optimal Culpability in Research Teams**\*

Erika Dömötör<sup>†</sup>

Zacharias Maniadis<sup>‡</sup>

Nikolas Tsakas<sup>§</sup>

University of Cyprus

University of Cyprus

University of Cyprus

University of Southampton

September 24, 2024

#### Abstract

Recent scandals in science have brought attention to the problem of detecting fraud and attributing punishment in the context of increasingly large research teams. We examine the problem theoretically and consider the socially optimal scheme for assigning culpability. We consider the simplest possible environment with two scientists, only one of whom is capable of committing fraud. Our theoretical analysis shows that a regime of group accountability that incentivises researchers to monitor other members of the group achieves the best social outcomes. Given this regime, the model yields the counter-intuitive prescription that punishing non-culpable members of the team for participating in a fraudulent project is the most promising tool for increasing the fraction of research that is honest.

Key words: Scientific Misconduct, Punishment, Monitoring, Research Teams JEL Codes: C72, D83, K30

<sup>\*</sup>This project has received funding from the European Union's Horizon 2020 Research and Innovation Programme under Grant Agreement number: 857636 – SInnoPSis – H2020-WIDESPREAD-2018-2020/H2020-WIDESPREAD-2018-04. This project has also received funding from the European Union's Horizon Europe research and innovation programme under grant agreement No [101079196] (Twinning for Excellence in Management and Economics of Research and Innovation [TWIN4MERIT]). We are grateful to Thomas Gall, Aristotelis Boukouras, Lorenzo Zirulia, David K. Levine and seminar participants at the University of Southampton.

<sup>&</sup>lt;sup>†</sup>Corresponding author. Address: SInnoPSis Research Group, University of Cyprus, 1 Panepistimiou Avenue, 2109 Aglantzia, Cyprus, email: domotor.erika@ucy.ac.cy

<sup>&</sup>lt;sup>‡</sup>SInnoPSis Research Group, University of Cyprus, 1 Panepistimiou Avenue, 2109 Aglantzia, Cyprus, email: maniadis.zacharias@ucy.ac.cy

<sup>&</sup>lt;sup>§</sup>Department of Economics, University of Cyprus, P.O. Box 20537, Nicosia 1678, Cyprus. email: tsakas.nikolaos@ucy.ac.cy

# 1 Introduction

The problem of scientific misconduct has recently attracted great attention, with the resignation of leading scientists who held top administrative positions at renowned institutions, following serious concerns about their research practices. In behavioural sciences, serious controversies surround allegations of systematic fraud and misconduct by world-famous researchers (Thorp, 2023). Many interventions have been proposed to tackle such problems, but their rigorous evaluation is difficult, due to the complexity of social systems for which Randomized Control Trials may be infeasible. Still, mathematical modelling of scientific institutions using game theory is a very promising avenue for simulating the effects of even such complex interventions (Gall et al., 2017). The issue of how to detect and deter misconduct is intimately related to the fact that research production takes place within teams, rather than individually. Accordingly, in designing scientific institutions robust to misconduct, one needs to combine the careful study of individual incentives with an understanding of the structure of scientific teams and its implications for individual incentives.

The importance of studying incentives for misconduct within a team setting is amplified by another important development in science: the constantly increasing average size of research teams, where larger teams are associated with more impact (Larivière et al., 2015). For instance, Poldrack et al. (2017) note that in genetics, neuroscience and other areas, larger consortia are needed to address the problem of low power of individual studies. This tendency may be beneficial for discovery, but it also raises new problems. Hall et al. (2018) argue that "... rapid increases in the demand for scientific collaborations have outpaced changes in the factors needed to support teams in science, such as institutional structures and policies, scientific culture, and funding opportunities." The issues of authorship, credit and accountability, and the consequences of an increasing average size of research teams for the prevalence of misconduct come to the forefront.

The problem of accountability in science relates to the ancient problem of attributing punishment for misconduct when it is known that the perpetrator belongs to a social group, but their identity is costly to reveal. The question is whether the group should have some sort of joint responsibility, and this issue has been related to law enforcement since antiquity (Miceli and Segerson, 2007). The main ethical concern is whether punishment of the innocent is acceptable as means to induce internal monitoring by other team members. Should collaborators of research fabricators be held accountable? Would that decrease fraudulent research? What will happen to the production of honest research?

In terms of actual policy, the scientific community is considering different approaches in assigning group responsibility in science (Helgesson and Eriksson, 2018; Hussinger and Pellens, 2019). Overall, the main approaches can be categorized as follows:

- A group accountability regime, where everyone in the team is responsible for any misconduct.
- A partial group accountability regime, where anyone in the team is responsible as long as they knew about the misconduct.
- Individual accountability, where only the persons committing misconduct are responsible.
- A guarantor regime, where an overall coordinator of the project is accountable for all aspects of a research paper.

Our main focus will be to study the optimal regime for the allocation of accountability in the case of wrongdoing in scientific projects, utilizing the categories of Helgesson and Eriksson (2018) and Hussinger and Pellens (2019). Our model focuses on an environment with two scientists (for tractability) who are asymmetric. Scientist 1 may cheat to improve the chances of publication of the final article, whereas Scientist 2 may choose to monitor Scientist 1 and to abort the project in the case of suspicious evidence appearing.<sup>1</sup> Our objective is to study the rate at which the research is conducted under the different regimes, and the degree to which this research is fraudulent.

For each regime that we study we shall assume that the enforcer pre-commits to a certain type of punishment, regardless of how unfair the punishment may look after the fact. Such proportionality considerations can be addressed at the design stage, but given an established regime, we assume

<sup>&</sup>lt;sup>1</sup>This asymmetry stems from the fact that Scientist 1 is interpreted as a junior researcher and Scientist 2 as a senior one. This assumption corresponds to many environments of scientific collaboration. For instance, there is a current tendency to use a contributorship model, where the task for each team member is declared in the scientific article. In this context, junior scientists are usually assigned tasks in which scientific fraud is possible (running experiments, data management, etc.) while senior scientists are assigned tasks for which fraud is not an issue (conceiving the study, revising the paper, etc). Accordingly, our setting applies to a general set of environments. Of course, in a different setting, where all researchers have the opportunity to cheat as well as to monitor others, the players would also face coordination issues with regards to their levels of cheating, which could create interesting equilibria. Modelling such an environment would require capturing asymmetric incentives and levels of power between coauthors.

that there is commitment. In addition, we shall assume that the punishment of collaborators principally depends on the established formal policies and regimes.<sup>2</sup>

Our model yields interesting results on the social benefits of alternative institutional and policy schemes. In cases of scientific misconduct with asymmetric researchers, the problem of deterrence arises as long as society's maximum penalty to fraudsters is bounded (e.g. career termination). In the absence of uncertainty about the identity of the perpetrator, our model reveals that the problem would be easily resolved as long as a sufficient punishment for the perpetrators was established. However, junior researchers with 'little to lose' will not be deterred by their potential punishment alone. The problem of deterrence remains as long as society is reluctant to punish further, e.g. by criminal prosecution.

Our model then shows that regimes that go beyond individual accountability may achieve better social outcomes. While the punishment may not be enough to deter potential fraud, group accountability may induce more senior collaborators to monitor and abandon projects that raise 'red flags'. Given the institution of group accountability, we show that the best policy for maximising the fraction of research that is honest is to maximally punish the non-culpable collaborator for participating in a fraudulent project.

Our results raise interesting questions regarding the tradeoff between the fairness and proportionality of punishment and the end outcomes for society. For societies unwilling to heavily (e.g. criminally) punish perpetrators but who would tolerate group punishment, a regime of group accountability is better in terms of its final outcomes. We should note that our model shows the advantages of group accountability even in the absence of the problem of identifying the perpetrator. If this problem exists, there may be further benefits by the group accountability regime.

The problem of scientific misconduct has attracted increasing attention in the last two decades across scientific disciplines (Fanelli, 2009; John, Loewenstein and Prelec, 2012), and several factors have been associated to it, such as academic culture as well as structural and psychological factors (Fanelli et al., 2015). The economics response to this problem (Lacetera and Zirulia, 2011) involves

<sup>&</sup>lt;sup>2</sup>Our approach also abstracts from the problem of identifying the culprit of a revealed misconduct – and associated 'detection costs' (as in Miceli and Segerson, 2007). Accordingly, we do not consider the fact that the accountability regime may increase incentives to provide information for identification of the perpetrator after misconduct has been confirmed. However, it is worth pointing out that in many domains there are very weak incentives to reveal a perpetrator (Vie, 2020).

the study of relevant incentive structures, and in particular the effects of policies for regulating and punishing such behaviors, in the spirit of the classical economics of crime literature (Becker, 1968). In fact, evidence indicates that such policies may have a primary role in explaining the existing evidence with regards to misbehavior (Fanelli et al., 2015).

Marx and Squintani (2009) provide a model of individual accountability in teams, assuming group punishment is not possible. Like our model, they assume that delegated monitoring does not have an informational advantage. Their model shows that delegated monitoring can be efficiencyenhancing – even in the absence of group penalty – as long as the monitoring requirement is included in worker's job assignments. Miceli and Segerson (2007) model the problem of determining under which conditions group punishment is preferred to individual punishment. Unlike our model, they consider the informational advantages of group punishment and explicitly model the social cost of punishing the innocent. Surprisingly, they find that if the goal of punishment is retribution (ex post fairness of punishment) group punishment is preferred to individual random punishment. Kiri et al. (2018) examine incentives to conduct verification activities in science, as opposed to producing novel research. Like our study, they find that (as long as any research is produced) a positive fraction of low-quality research characterizes any equilibrium of the game. As a result, society needs to tolerate some level of misbehaviour or low quality research as part of the normal functioning of science.

The microfinance literature has considered the problem of individual vs. joint liability for a long time. Ghatak and Guinnane (1999) review the literature, considering the advantages of group borrowing in the context of developing countries, where borrowers typically lack collateral. They illustrate how joint liability can affect the formation of teams, when different agents have different types (which could be interpreted in our model as the propensity to conduct misbehaviour). Joint liability can utilize local information and weed out the bad types. On the other hand, joint liability does not ameliorate moral hazard in this model, since players do not incorporate the effects of their own action on their partner's action. Similar to our model, joint liability improves welfare. This happens under certain conditions, associated with social sanctions that are effective enough or low monitoring costs. Karlan (2007) consider how social connections facilitate the monitoring and enforcement of loans in a setting of group liability. Moreover, De Quidt et al. (2016) show that, for high levels of social capital, group lending under individual liability can outperform joint liability.



Figure 1: Game Tree

## 2 The Model

Consider an environment with two scientists who may collaborate on a joint project. The two scientists differ with respect to the payoffs they may receive from the project. Moreover, one of them, scientist **S1** (he) may commit fraud to improve the outcome of the project. Scientist **S2** (she) is aware of this possibility, but cannot observe directly whether fraud has taken place. Nevertheless, **S2** has the option of monitoring **S1** by using a costly and imperfect monitoring mechanism. We model this using a dynamic game between **S1** and **S2**. The detailed game tree is presented in Figure 1.

First, S1 decides whether to commit fraud or be honest, i.e. the action set of player S1 is

 $A_1 := \{F, H\}$ . Accounting also for mixed strategies, the strategy of **S1** can be summarized by the probability  $p \in [0, 1]$  with which he commits fraud.

Subsequently, without observing the choice of **S1**, scientist **S2** decides whether to walk away from the project, proceed without monitoring, or proceed with monitoring. Thus, the action set of **S2** is  $A_2 := \{WA, NM, M\}$  and a (mixed) strategy of **S2** can be summarized by the pair of probabilities  $(q_1, q_2)$ , where  $q_1 := \mathbb{P}(WA)$  and  $q_2 := \mathbb{P}(NM)$ . Obviously, the remaining probability  $1 - q_1 - q_2 = \mathbb{P}(M)$ , while  $q_1, q_2 \ge 0$  and  $q_1 + q_2 \le 1$ .

If **S2** decides to monitor, then she also chooses the precision of monitoring,  $z \in (0, 1/2]$ . Namely, conditional on **S2** having chosen to monitor, **S2** performs a costly investigation. Monitoring precision z implies that the outcome of the investigation is correct with probability 1/2 + z. That is, **S2** observes an investigation outcome in  $\{f, h\}$  such that  $\mathbb{P}(f|F) = \mathbb{P}(h|H) = 1/2 + z$ . Monitoring is costly, with the cost of choosing monitoring precision z being  $c(z) = \frac{1}{2}z^2$ . After observing the signal, **S2** can decide the walk away from the project at this point (WA) or to proceed with the project (P).

Notice that, if S2 decides to monitor, but does so with a very low intensity, then the signal would be largely uninformative, thus no realization would be able to alter her choice. Thus, following the investigation's suggestion would be suboptimal for S2. However, as we show in detail in Section 3, the costly nature of monitoring implies that any strategy in which S2 monitors but does not follow the signal obtained by the investigation is strictly dominated. This allows us to consider a simplified version of the game, which shares the same equilibrium outcomes as the original game, while also being handy in the subsequent analysis. Namely, we consider that S2 commits to always comply with the recommendation of the investigation, which in our model is equivalent to committing to follow the signal she receives.<sup>3</sup> On an intuitive level, the reduced model in which S2 is willing to participate in a project that might end up being fraudulent, only as long as she can claim plausible deniability if fraud gets detected later on.

Once the players have made their choices, payoffs are realized. If the project is abandoned because S2 walks away – either immediately, or after observing the signal – then the players

<sup>&</sup>lt;sup>3</sup>Alternatively, this modeling choice regarding monitoring can also be seen as an economical way to capture the existence of a fixed cost in monitoring (i.e. independent of the monitoring intensity).

obtain the payoff of some outside option  $\overline{u}_1$  and  $\overline{u}_2$ , respectively. If the project is completed, then the expected payoffs for the two players depend on whether **S1** has committed fraud or not, on the regulatory regime, and in some cases on whether or not **S2** knew about the fraud. If **S1** has been honest, then the payoffs from completing the project are  $R_1$  and  $R_2$ , respectively. On the other hand, if he has committed fraud the outcomes are either  $\hat{R}_1$  and  $\hat{R}_2$  if the fraud is not caught, or  $R_1 - \tau_1$  and  $R_2 - \tau_2$  if it is caught by the external authority. We assume that the external authority works with a monitoring technology which identifies a fraud with  $\pi$  success rate, and never makes false accusations, i.e. never identifies honest research as fraudulent by mistake. Using these parameters, the expected payoffs from fraudulent research are  $\underline{R}_1 = (1 - \pi)\hat{R}_1 + \pi(R_1 - \tau_1)$ for **S1**, and  $\underline{R}_2 = (1 - \pi)\hat{R}_2 + \pi(R_2 - \tau_2)$  for **S2**.

For simplicity, we assume that  $\overline{u}_1$ ,  $R_1$ , and  $\underline{R}_1$  are all different from each other, and  $\overline{u}_2$ ,  $R_2$ , and  $\underline{R}_2$  are all different from each other. Furthermore, we assume that fraudulent research has the potential to make both researchers better off in case it is not caught, that is  $\hat{R}_1 > R_1$  and  $\hat{R}_2 > R_2$ . We also set lower bounds of the punishment at  $\tau_1 \ge R_1$  and  $\tau_2 \ge R_2$ , assuming that when fraud is caught the value of the research is lost, even if there is no further sanction.

Under some regulatory regimes it is possible that the punishment for **S2** if there is a fraud depends on whether or not she knew about it. To cover these cases we enrich our notation such that the outcome for **S2** from fraudulent research is  $\underline{R}_2^{NM}$  if she did not attempt to monitor at all,  $\underline{R}_2^{Mh}$  if she monitored but the signal suggested honest research, and  $\underline{R}_2^{Mf}$  if she monitored and got signal of fraudulent research, that is, if she knew about the fraud.

## **3** Equilibrium Analysis

In this section we specify the parameters in case of each regulatory regime and solve for the equilibrium.<sup>4</sup> We focus on the cases where  $\overline{u}_1 < R_1 < \underline{R}_1$ , that is when **S1** has an incentive to commit fraud even under the threat of punishment.<sup>5</sup> We depart from the individual accountability regime

<sup>&</sup>lt;sup>4</sup>The natural equilibrium notion to consider here is Perfect Bayesian Equilibrium. Yet, in two of the regimes, we consider a simplified version of the game, which makes it essentially a simultaneous move game. To avoid unnecessary complications, we describe the equilibria referring only to the equilibrium strategy profile.

<sup>&</sup>lt;sup>5</sup>In case it was possible for the social planner to set such high punishment for **S1** that  $\underline{R}_1 < R_1$ , then the analysis would trivially lead to the first-best outcome of no fraud. The threshold for this high enough punishment would be  $\tau_1 > \frac{1-\pi}{\pi}(\hat{R}_1 - R_1)$ . However, we assume that such high punishment is not possible, for example, because the reward from a fraudulent research  $\hat{R}_1$  is too high, or the monitoring accuracy  $\pi$  is too low.

and show that under the former assumption the social planner has no policy tool to prevent fraud. Then, we continue with the partial group accountability regime, arguably a morally defensible option. However, we can quickly conclude that it is still not compatible with economic incentives, and the threat of punishing **S2** in this scheme does not increase internal monitoring or discourage fraudulent research. Therefore, we explore alternative regulatory regimes; the group accountability and guarantor regimes. Throughout the analysis we will assume (without loss of generality) that the outside option is normalized at  $\overline{u}_1 = \overline{u}_2 = 0$ .

#### 3.1 Individual Accountability Regime

Under the individual accountability regime we assume that the punishment for **S2** is merely the loss of the value of the project:  $\tau_2^{NM} = \tau_2^{Mh} = \tau_2^{Mf} = R_2$ , therefore  $\underline{R}_2 = \underline{R}_2^{NM} = \underline{R}_2^{Mh} = \underline{R}_2^{Mf} = (1 - \pi)\hat{R}_2$ . Since  $\underline{R}_2 > 0$ , Walk Away is strictly dominated by Not Monitor. Similarly, Monitor is strictly dominated by Not Monitor, because in case of fraudulent research  $\underline{R}_2$  is larger than any possible outcome after monitoring, either -c(z) or  $\underline{R}_2 - c(z)$ , and in case of honest research  $R_2$ is larger than any outcome after monitoring, either -c(z) or  $R_2 - c(z)$ . Knowing that **S2** will not monitor, **S1** chooses to commit fraud with probability 1. Therefore, the unique equilibrium is (Fraud, Not Monitor), which is presumably the worst of all outcomes from the social planner's perspective.<sup>6</sup>

### 3.2 Partial Group Accountability Regime

In the partial group accountability regime the scientist who commits fraud is responsible for his actions, and additionally, the other scientist who is in the position to monitor is also culpable, but only if she knew about the fraud. In terms of our model, we assume that the punishment for **S1** is larger than the minimal:  $\tau_1 > R_1$ . **S2** faces a larger-than-minimal punishment only in case she monitors and gets a signal about the fraud:  $\tau_2^{Mf} > R_2$ . If she does not monitor, or if she monitors but does not detect the fraud, she merely loses the value of the project:  $\tau_2^{NM} = \tau_2^{Mh} = R_2$ . Her

 $<sup>^{6}</sup>$ We postpone the formal introduction of the objective of the social planner until Section 4, where we assume that the planner's aim is to minimize fraudulent research and maximize honest research.

expected payoff from the latter case is  $\underline{R}_2^{NM} = \underline{R}_2^{Mh} = (1 - \pi)\hat{R}_2 > 0$ , which also means that  $\underline{R}_2^{NM} = \underline{R}_2^{Mh} > \underline{R}_2^{Mf}$ .

Since  $\underline{R}_2^{NM} > 0$ , Walk Away is strictly dominated by Not Monitor. Similarly, Monitor is strictly dominated by Not Monitor, because in case of fraudulent research  $\underline{R}_2^{NM}$  is larger than any possible outcome after monitoring: -c(z),  $\underline{R}_2^{Mf} - c(z)$ , or  $\underline{R}_2^{Mh} - c(z)$ , and in case of honest research  $R_2$  is larger than any outcome after monitoring, either -c(z) or  $R_2 - c(z)$ . Knowing that **S2** will not monitor, **S1** chooses to commit fraud with probability 1, hence the unique equilibrium is (Fraud, Not Monitor), just like in the individual accountability regime.

We conclude that the policy maker does not gain from partial group accountability compared to individual accountability. This result is driven by the fact that S2 can get away from the punishment by simply not trying to monitor. Therefore, we need to investigate further regimes in which S2, the scientist who is capable of internal monitoring, is unconditionally culpable.<sup>7</sup>

### 3.3 Group Accountability Regime

In this regime, we work with the general case where both **S1** and **S2** are culpable, regardless of **S2**'s choice of monitoring and the observed signal. However, the sizes of the punishments for the two scientists are not necessarily equal, we only specify that  $\tau_1 > R_1$  and  $\tau_2^{NM} = \tau_2^{Mh} = \tau_2^{Mf} > R_2$  and therefore  $\underline{R}_2^{NM} = \underline{R}_2^{Mh} = \underline{R}_2^{Mf} = \underline{R}_2$ . Furthermore, we assume that  $\underline{R}_2 < 0 < R_2$ .<sup>8</sup> The equilibrium in this case depends on the exact parameter choices.

#### 3.3.1 Preliminaries

Before proceeding with the detailed analysis, we shall make a few observations. First, we show that any strategy in which **S2** does not follow her signal is strictly dominated:

<sup>&</sup>lt;sup>7</sup>Assuming an alternative environment where there is imperfect capacity of the external authority to identify cases where monitoring has taken place, the incentives of S2 to monitor would be stronger than before, yet still not sufficient to improve over *Not Monitor*.

<sup>&</sup>lt;sup>8</sup>If this assumption is violated **S2** never has an incentive to monitor. If  $\underline{R}_2 > 0$  we would get a trivial solution, similar to the individual and partial group accountability regime, where **S2** never monitors and **S1** always commits fraud. Thus, the group accountability regime is only effective in promoting monitoring and preventing fraud when the punishment is large enough to make **S2** worse off staying in a fraudulent project rather than walking away.



Figure 2: Reduced Game Tree under Group Accountability

- 1. The pooling strategy Walk Away after either f or h signal which yields payoff -c(z) is strictly dominated by Walk Away without monitoring which yields 0.
- 2. The pooling strategy *Proceed* after either f or h signal which yields payoff  $p\underline{R}_2 + (1-p)R_2 c(z)$  is strictly dominated by *Not Monitoring* which yields  $p\underline{R}_2 + (1-p)R_2$ .
- 3. The separating strategy Walk Away after h and Proceed after f which yields expected payoff  $p(1/2+z)\underline{R}_2 + (1-p)(1/2-z)R_2 c(z)$  is dominated by the separating strategy Proceed after h and Walk Away after f which yields expected payoff  $p(1/2-z)\underline{R}_2 + (1-p)(1/2+z)R_2 c(z)$ .

Since we eliminated the strategies where **S2** does not follow her signal, we can reduce the game as shown in Figure 2.

Second, the reduced game has a single subgame, where the action set of S2 regarding the precision of monitoring is a whole interval. Yet, the problem can be simplified substantially.

Observe that the expected utility of **S2** conditional on having chosen to monitor (M) and on believing that **S1** has committed fraud with probability p is given by the following expression:

$$EU_2(z|p,M) = p\left[\left(\frac{1}{2} - z\right)\underline{R}_2\right] + (1-p)\left[\left(\frac{1}{2} + z\right)R_2\right] - \frac{1}{2}z^2$$
(1)

It is straightforward to observe that the above function is strictly concave in z for all  $z \in (0, 1/2]$ , and thus it attains a unique maximum within that range, which may not be interior.

Therefore, we can simplify the game by substituting, for each belief p, the optimal monitoring intensity  $z^*(p)$  to the expected utility of **S2** if she chooses to monitor.<sup>9</sup> Namely, for scientist **S2**, we obtain:

$$EU_{2}(M|p) = p\left[\left(\frac{1}{2} - z^{*}(p)\right)\underline{R}_{2}\right] + (1-p)\left[\left(\frac{1}{2} + z^{*}(p)\right)R_{2}\right] - \frac{1}{2}\left[z^{*}(p)\right]^{2}$$

$$EU_{2}(NM|p) = p\underline{R}_{2} + (1-p)R_{2}$$

$$EU_{2}(WA|p) = 0$$
(2)

whereas for S1 we obtain:

$$EU_{1}(H|q_{1},q_{2},z) = q_{2}R_{1} + (1-q_{1}-q_{2})\left(\frac{1}{2}+z\right)R_{1}$$

$$EU_{1}(F|q_{1},q_{2},z) = q_{2}\underline{R}_{1} + (1-q_{1}-q_{2})\left(\frac{1}{2}-z\right)\underline{R}_{1}$$
(3)

Note that the reduced game can be seen as a simultaneous move game, and an equilibrium would then be characterized by a triplet  $(p^*, q_1^*, q_2^*)$  and by  $z^*(p^*)$ .<sup>10</sup> We also define the following quantities: (i)  $a := -\underline{R}_2 > 0$ , which is the loss **S2** incurs by remaining involved in a fraudulent project, (ii)  $b := R_2 > 0$ , which is the gain **S2** enjoys by remaining involved in an honest project, and (iii)  $k = \frac{1}{2} \frac{\underline{R}_1 - R_1}{\underline{R}_1 + R_1}$ , which is a parameter that measures **S1**'s relative gain from committing fraud.

<sup>&</sup>lt;sup>9</sup>Of course, this expected utility is no longer going to be linear in p.

<sup>&</sup>lt;sup>10</sup>Formally, the notion of equilibrium in the original game to be used would be Perfect Bayesian Equilibrium, where the equilibrium beliefs of **S2** in both information sets would be derived from  $p^*$  using Bayes rule. Therefore, the posterior belief after signal f is  $p^f = \frac{p^*(1/2+z)}{p^*(1/2+z)+(1-p^*)(1/2-z)}$ , and after signal h is  $p^h = \frac{p^*(1/2-z)}{p^*(1/2-z)+(1-p^*)(1/2+z)}$ . In the final stage, the action of **S2** would be WA after signal f and P after signal h.

Note that  $k \in [k_{min}, k_{max}] \subset (0, 1/2)$ , where  $k_{min}$  is defined by the maximum punishment  $\tau_1^{max}$  the social planner could impose on **S1**, which, by our assumption, is still not enough to fully eliminate his incentive to commit fraud:  $\underline{R}_1^{min} = (1-\pi)\hat{R}_1 + \pi(R_1 - \tau_1^{max}) > R_1$ . On the other hand, the minimum punishment for **S1** is to lose the value of the project, therefore  $k_{max}$  takes place when  $\tau_1 = \tau_1^{min} = R_1$ . We also assume that  $a \in [a_{min}, a_{max}] \subset (0, 1/2)$  and  $b \in [b_{min}, b_{max}] \subset (0, 1/2)$ .

#### 3.3.2 Best Responses and Equilibrium Characterization

We start by calculating the optimal monitoring precision S2 chooses upon deciding to monitor.

$$\frac{\partial EU_2}{\partial z} = pa + (1-p)b - z \Rightarrow z^*(p) = pa + (1-p)b$$

Notice that for  $a, b \in (0, 1/2), z^*(p) \in (0, 1/2)$  for every  $p \in [0, 1]$ .

Given the optimal precision of monitoring calculated above, the expressions in (2) regarding **S2**'s expected utility can be rewritten as follows:

$$EU_{2}(M|p) = \frac{1}{2} \left[ -pa + (1-p)b \right] + \frac{1}{2} \left[ pa + (1-p)b \right]^{2}$$
  

$$EU_{2}(NM|p) = -pa + (1-p)b$$
  

$$EU_{2}(WA|p) = 0$$
(4)

We can now calculate the best-responses of **S2** for all p. The detailed calculations can be found in the Appendix. An important observation is that **S2**'s best response can be split into three areas depending on p. Namely, **S2** stays in without monitoring (chooses NM) when p is small, stays in and monitors (chooses M) when p is moderate, and walks away (chooses WA) when pis high. These three regions are determined by two thresholds  $p_1^*$  and  $p_2^*$ , as shown on Figure 3. Both thresholds can be explicitly calculated (see Expression 6 in the Appendix) and are shown to be between 0 and 1 and such that  $p_1^* < p_2^*$ . This also means that randomization may occur only between NM and M for  $p = p_1^*$ , or between M and WA for  $p = p_2^*$ .

Let us now calculate the best response of **S1** for each  $(q_1, q_2)$  and z. Using the quantity  $k = \frac{1}{2} \frac{R_1 - R_1}{R_1 + R_1}$ , we find that the optimal decision of **S1** depends on whether the incentives to commit



Figure 3: Strategy of S2

fraud are sufficiently high given  $(q_1, q_2)$  and z, which depends on the sign of the expression  $2q_2k + (1 - q_1 - q_2)(k - z)$ .

Having described the best responses of the two players, we can proceed to obtain the equilibria of the game. Proposition 1 presents all equilibria of the game, with the respective parameter conditions in which each equilibrium exists. The parameter conditions are not mutually exclusive, thus there are parameter values for which the game has multiple equilibria.

**Proposition 1** For  $a \in [a_{min}, a_{max}] \subset (0, 1/2)$ ,  $b \in [b_{min}, b_{max}] \subset (0, 1/2)$ , and  $k \in [k_{min}, k_{max}] \subset (0, 1/2)$  the game has the following types of equilibria:

$$\begin{array}{l} \textbf{[1]} (\forall a, b, k): p^* \in [p_2^*, 1], \ (q_1^*, q_2^*) = (1, 0), \ z^* = p^*a + (1 - p^*)b. \\ \textbf{[2]} (a \neq b \ and \ \frac{k-b}{a-b} \in [p_1^*, p_2^*]): p^* = \frac{k-b}{a-b}, \ (q_1^*, q_2^*) = (0, 0), \ z^* = k. \\ \textbf{[3a]} (a \neq b \ and \ \frac{k-b}{a-b} = p_2^*): p^* = p_2^*, \ (q_1^*, q_2^*) = (\hat{q}, 0) \ for \ \hat{q} \in (0, 1), \ z^* = k. \\ \textbf{[3b]} (a = b = k): p^* = p_2^*, \ (q_1^*, q_2^*) = (\hat{q}, 0) \ for \ \hat{q} \in [0, 1), \ z^* = k. \\ \textbf{[4a]} (a > b \ and \ \frac{k-b}{a-b} \leq p_1^*): p^* = p_1^*, \ (q_1^*, q_2^*) = \left(0, \frac{p_1^*a+(1-p_1^*)b-k}{p_1^*a+(1-p_1^*)b+k}\right), \ z^* = p_1^*a + (1 - p_1^*)b. \\ \textbf{[4b]} (a < b \ and \ \frac{k-b}{a-b} \geq p_1^*): p^* = p_1^*, \ (q_1^*, q_2^*) = \left(0, \frac{p_1^*a+(1-p_1^*)b-k}{p_1^*a+(1-p_1^*)b+k}\right), \ z^* = p_1^*a + (1 - p_1^*)b. \\ \textbf{[4c]} (a = b \geq k): p^* = p_1^*, \ (q_1^*, q_2^*) = \left(0, \frac{a-k}{a+k}\right), \ z^* = a. \end{array}$$

From this result, we can get some useful observations. First of all, irrespective of the parameter values, there is always an equilibrium in which S2 walks away (eqm. 1), because S1 commits fraud with sufficiently high probability. This equilibrium captures the case in which there is lack of trust between the two players. There is also an equilibrium in which S2 mixes between monitoring and walking away (eqm. 3), but this can be sustained only for very specific combinations of parameter values, thus it is of limited use.

In addition to the above, there are two equilibria that involve monitoring by S2. There is one equilibrium (eqm. 2) in which S1 commits fraud with a probability between  $p_1^*$  and  $p_2^*$  and for which S2 always monitors. There is also one equilibrium (eqm. 4) in which S1 commits fraud with a probability exactly equal to  $p_1^*$  and for which S2 mixes between monitoring and non-monitoring. Both of these types of equilibria are preferred by both players compared to the equilibrium in which S2 always walks away, which yields payoff 0 to both players. This happens because a payoff equal to zero could always be achieved for both players with certainty if S2 were to walk away. The fact that S2 does not do so in equilibrium implies that she shall expect a higher payoff.

What about the other two sensible equilibria (the one with certain monitoring and the one with mixing between monitoring and non-monitoring)? Interestingly, in the parameter conditions for which both of them exist, the equilibrium with probabilistic monitoring yields higher payoffs for both players.

**Corollary 1** For a > b and  $\frac{k-b}{a-b} \in (p_1^*, p_2^*)$ , the game has three equilibria: **1**, **2**, and **4**. Within this range, among the three equilibria, eqm **4** yields the highest expected payoff to both players.

For S2 this happens because the equilibrium with certain monitoring is associated with a higher probability of fraud, which is detrimental for S2's payoffs, given that  $EU_2(M|p)$  decreases in p. For S1 the result is driven by two facts. On the one hand, fraudulent projects (that are likely to be stopped if monitored) are monitored less often when monitoring is probabilistic. On the other hand, honest projects are less often abandoned because of an incorrect signal, because equilibrium precision  $z^*$  is higher when monitoring is probabilistic.

To summarise, eqm. 1 involves the uninteresting case where S2 always walks away and mutual trust between players fails. On the other hand, eqm. 3 obtains for very specific parameters and is

therefore of limited relevance. Finally, eqm. 4 dominates eqm. 2 in terms of efficiency. Therefore, given its welfare advantage, from now on we will focus on the analysis of eqm. 4, which involves probabilistic monitoring.

### 3.4 Guarantor Accountability Regime

In the guarantor accountability regime only **S2** is held responsible for any fraud (which only **S1** can commit), and her punishment is  $\tau_2^{NM} = \tau_2^{Mf} = \tau_2^{Mh} > R_2$ , while the punishment for **S1** is  $\tau_1 = R_1$ . This case can be seen as a special case of the group accountability regime. The conditions  $\underline{R}_1 \ge R_1 \ge 0$  and  $\underline{R}_2 \le 0 \le R_2$  from the main part are met for a range of different parameters. Therefore, the equilibria found in Proposition 1 still hold, and the preferred equilibrium is eqm. **4** by Corollary 1. The only difference compared to the analysis of the group accountability regime is that k is fixed at its maximum value in this case, which will be a relevant distinction in comparative statics and policy application.

## 4 The Quantity and Quality of Produced Research

In this section, we turn our attention to the impact of scientific fraud on research from a social point of view. In line with the features of our model, we focus on the effect of the parameter values on research output, namely on the expected quantity of honest and fraudulent research. The parameters of interest are a, b, and k, among which a and k are associated with the punishment of fraudulent research, thus are affected by the accountability regime.

In general, and in absence of any constraints, a planner would like to design these incentives so as to both maximize the production of honest research (henceforth called HR), and, at the same time, minimize the production of fraudulent research (henceforth called FR). Yet, in the presence of constraints, these two goals may require conflicting choices, thus generating trade-offs in the planner's policy settings. To capture these trade-offs, we consider that the planner intends to maximize the difference between the two quantities, i.e. maximize HR-FR. This assumption implies that the planner cares equally about the two objectives. Nevertheless, it will become apparent that the result would be qualitatively similar for more general functions that are linear combinations of the two quantities.

### 4.1 Comparative Statics

To provide a more general understanding, we begin this section by providing some comparative statics, regarding the effects of the parameters of the model on each measure of research output. Throughout this section, and based on the results of the previous section, we focus on the mixed  $Monitor/Not\ Monitor$  equilibrium of the Group Accountability Regime (Equilibrium 4 in Proposition 1). We are interested both in the level of total research output, as well as in the level of fraudulent research that is eventually published. The following expressions quantify equilibrium levels of total research (TR), fraudulent research (FR), and honest research (HR) as functions of a, b, and k, where a and b appear implicitly in the expressions through  $p_1^*$  and  $z^*$ :

$$TR = q_2^* + (1 - q_2^*) \left( p_1^* \left( \frac{1}{2} - z \right) + (1 - p_1^*) \left( \frac{1}{2} + z \right) \right) = z^* \frac{1 + 2k(1 - 2p_1^*)}{z^* + k}$$
$$HR = (1 - p_1^*) \left( q_2^* + (1 - q_2^*) \left( \frac{1}{2} + z^* \right) \right) = z^* (1 - p_1^*) \frac{1 + 2k}{z^* + k}$$
$$FR = p_1^* \frac{z^* (1 - 2k)}{z^* + k}$$

**Proposition 2** In the space of parameters (a, b, k) where the mixed equilibrium of the Group Accountability Regime (Equilibrium 4 of Proposition 1) exists, the following results hold:

- 1. For all a, b, Total Research (TR), Honest Research (HR), and Fraudulent Research (FR) decrease in k.
- 2. For all a, b, the ratio HR/TR increases in k.
- 3. For all b, k, TR and HR increase in a.
- 4. For all b, (i) for all  $k > \frac{\sqrt{29}-3}{10}$ , FR decreases in a, whereas (ii) for all  $k < \frac{\sqrt{29}-3}{10}$ , there exists  $\hat{a}(b) \in (0, b)$  such that FR increases in a for  $a < \hat{a}$  and decreases for  $a > \hat{a}$ .
- 5. For all b, k, the ratio HR/TR increases in a.
- 6. For all a, k, FR increases in b.
- 7. For all a, k, the ratio HR/TR decreases in b.

The comparative statics presented in Proposition 2 provide interesting insights regarding the impact of the parameters on the equilibrium quantities of interest. First, the measure of the incentive to commit fraud denoted by k is inversely related to the punishment level of **S1**. Importantly, an increased incentive for **S1** to commit fraud does not necessarily have a detrimental effect. As long as k remains sufficiently small compared to a and b (so as not to preclude the existence of the mixed equilibrium), an increase in the value of k is internalized by **S2**, in the form of a higher equilibrium probability of monitoring. This leads to a decrease in both FR and HR, but fraudulent research decreases proportionally more. Nevertheless, if the incentives to commit fraud become too high, **S2** cannot counterbalance them through monitoring, which breaks down the mixed equilibrium and leads her to walk away from the project.

Second, there appears to be a more clearly detrimental effect from an increase in the value of b, which captures the value of a legitimate project to **S2**. An increase in b leads to more fraudulent research in absolute terms and to a lower share of honest research to total research. This happens because an increase in b leads both to less monitoring (albeit more accurate) and to a higher equilibrium probability of **S1** committing fraud.

On the contrary, a decrease in the expected payoff of S2 from being involved in a fraudulent project (i.e. an increase of a) has a clearly positive effect. It leads to higher levels of total research, which arise essentially exclusively from an increase in honest research, because more often than not fraudulent research decreases. This is due to the fact that a higher punishment leads to lower equilibrium levels of fraud and to a higher monitoring accuracy (albeit monitoring occurs less often).

Having said the above, it is interesting to observe that a and b have an opposite effect in the equilibrium probability of fraud  $(p_1^*)$ , despite the fact that both affect monitoring accuracy positively  $(z^*)$  and monitoring frequency negatively  $(1 - q_2^*)$ .

Overall, it seems that a promising tool to achieve the planner's objectives is to make the participation in a fraudulent project less appealing for **S2**. Note that the payoffs from participating to a fraudulent project implicitly depend on the probability of being detected by an external monitor (not a collaborator), and the accompanying punishment. Hence, this variable can to some degree be affected by norms and policies in the given scientific domain.

### 4.2 Optimizing Quantity and Composition of Research

We now turn our attention to the aggregate production of research. Namely, we would like to understand which values of the parameters would allow the social planner, who cares about maximizing honest research and minimizing fraudulent research, to achieve their goal most effectively.

Having said the above, the first-best outcome for such a social planner would be to implement an accountability policy (i.e. values of a and k) for which **S1** never commits fraud and **S2** always participates in the project without monitoring. Yet, note that this is not possible within the parameter range we have been studying, as it becomes apparent from Propositions 1 and 2. In fact, this could only be achieved if the incentives of **S1** to commit fraud were completely eliminated, i.e. if  $R_1 > \underline{R1}$  or equivalently k < 0. Thus, if a planner would be able to eliminate those incentives, then the problem under study would be solved. In what follows, we assume that the planner is not able to eliminate the incentives of **S1** to commit fraud.<sup>11</sup>

Importantly, also, as long as the incentives of **S1** to commit fraud (k > 0) are not completely eliminated, the only equilibrium in which no fraudulent research takes place is the one in which **S2** always walks away, so that in general no research is conducted. Therefore, the social planner is somewhat bound to tolerate some level of fraudulent research.

Given this observation, as described above, we consider the optimal choice of a planner who intends to maximize the difference between honest and fraudulent research, i.e. HR-FR. Note that the form of the objective function induces a tradeoff both for k (given that both HR and FR decrease in k) and for a, given that for some values of k and b both HR and FR increase in a.

**Proposition 3** In the space of parameters (a, b, k) where the mixed equilibrium of the Group Accountability Regime (Equilibrium 4 of Proposition 1) exists, there exists some  $\hat{b} \in (0, 1/2)$  such that HR-FR is maximized for

•  $(a,k) = (a_{max}, k_{min})$  if  $b \le \hat{b}$ , and

<sup>&</sup>lt;sup>11</sup>Intuitively, this means that in our environment there is a non-trivial problem only as long as society cannot impose a very large punishment to junior researchers who commit fraud. This is a reasonable assumption given the fact that reputation costs of juniors are low, and the termination of a career is the maximum penalty. Interestingly, even if fraudsters are punished under civil law, since juniors have not accumulated high earnings from fraud, the level of punishment seems limited. The question remains whether criminal law should be used more frequently in cases of scientific misconduct.

•  $(a,k) = (a_{max}, k_{max})$  if  $b > \hat{b}$ .

The result stems from the observation that HR - FR is never maximized for intermediate levels of k, while for all k the objective function is increasing in a.

The result yields an important takeaway. Namely, it is always beneficial for the planner to induce the maximum allowed punishment to the scientist who is responsible for monitoring, i.e. **S2**. By doing so, the planner leads **S2** to increase her effort of preventing fraud, while still participating in the project, which in turn leads **S1** to commit fraud with a smaller probability. The optimal punishment for **S1** is less intuitive and less clear. In fact, in some cases, it may be optimal to eliminate completely the punishment for **S1**. Arguably, such a policy may raise ethical concerns, as it leaves unpunished the person who is primarily responsible for fraud.

Note that in our terminology of accountability regimes, Proposition 3 suggests that for low rewards of S2 (low b's) the group accountability regime is optimal with the maximum possible punishment for both scientists. However, when the rewards are high enough for S2, the guarantor regime yields the best outcome. Intuitively, in our setting imposing the highest punishment on the person who is in position to monitor is always necessary, and it is also sufficient if this person has high enough stakes in the game; punishing the perpetrator in this case does not increase the social planner's objective further.

At this point, it is important to notice that the choice to maximize HR - FR implicitly assumes that the social planner cares equally about the two quantities. Yet, this need not always be the case. One could consider a more general version of this objective function of the form  $\lambda HR - (1 - \lambda)FR$ . In light of the above result, the outcome would be qualitatively similar to the previous one. It would still be optimal to target a = 1/2, and optimal k would depend on b, but also on  $\lambda$ , with high  $\lambda$  favoring low k and low  $\lambda$  favoring a high k.

Furthermore, so far we have assumed that the social planner faces no restriction on how to allocate culpability, other than the restrictions on the parameters' range imposed by the model. However, there may be other restrictions, constraining further the options of the planner, for instance due to the monitoring mechanisms available to the social planner herself. For example, creating mechanisms able to capture fraudulent research with very high probability would require a large amount of resources, which may be both expensive and inefficient. On the other hand, imposing really severe punishments to offenders may be inconsistent with more general legal practices. For instance, how severe a case of academic fraud would have to be for a court to justify imposing jail time to the offender? Such additional restrictions would likely affect the policy implications suggested by our analysis and they could be incorporated in the model. Our model serves as a benchmark for future research and further analysis.

# 5 Conclusion

The problem of scientific misbehaviour is attracting increasing attention given the alleged credibility crises in several disciplines, and is exacerbated by the increasing size of scientific teams. Several sets of rules have been suggested for the allocation of responsibility in case of scientific fraud. The literature has identified four different options: a group accountability regime, an individual accountability regime, a guarantor regime, and a partial group accountability regime. In terms of a purely ethical approach, the best solution is punishing the individuals who committed fraud, or at most those who knew about the fraud. However, our analysis shows that society may have to bear with a trade-off between ethical principles and efficiency of incentives. The reason is that, generally, the best outcomes for society are attained when the whole scientific team is responsible for scientific fraud (group accountability regime). This is driven by the fact that this scheme provides incentives for internal monitoring of other team members. We further show that the policymaker has to accept some level of fraud in order to facilitate any amount of joint projects.

It should be noted that only if the punishment level of the perpetrators can be arbitrarily high, then the first best outcome, all research being honest, is achievable and it does not require group accountability. However, social norms may prevent the policy-maker from setting such high individual punishment. We have assumed throughout that such limitations exist and necessitate alternative schemes of allocating accountability.

Our findings in the main analysis suggest that as long as S1 has an incentive to commit fraud, but S2 has no incentive to be involved in a fraudulent project, the game theoretic equilibrium is either not to start any joint research project, or to let S1 to have a mixture or fraudulent and honest research projects, while S2 monitors randomly a fraction of these projects, with an imperfect monitoring strategy. While the first equilibrium is clearly undesirable, in the second one the output consists of some honest research, some fraudulent research that was not caught during the internal and external monitoring processes, and some honest research which was abandoned by mistake due to the imperfection of the internal monitoring process. The evaluation of the combination of these various outcomes depends on the policy maker's objectives. However, the main alternative schemes (individual accountability and partial group accountability), achieve a clearly worse outcome, because they fail to provide incentives for internal monitoring. Accordingly, the team member who has the opportunity to cheat does so with probability 1.

Our results inform the discussion about optimal institutions in science. Our mathematical simulation allows us to analyse a complex institutional problem that cannot be examined with the experimental method. Our results suggest that internal monitoring in science can play a crucial role in limiting misbehaviour. It is worth emphasizing that we obtain this result without assuming that there is uncertainty regarding the identity of the perpetrator. In such an enhanced setting, the benefits of group accountability would likely be reinforced. Our analysis assumes only two researchers, whereas some of the key challenges involve increasingly large research teams. Accordingly, our approach would yield additional interesting insights if it was generalized to the case of multiple researchers. We leave such analysis for further work.

# Appendix

### **Calculations and Proofs**

**Best Response of S2:** We need to consider separately the cases where a = b and  $a \neq b$ .

Let us first consider the case a = b. In this case, we can simplify the expressions of the expected utilities as follows:  $EU_2(M|p) = \frac{(1-2p)a+a^2}{2}$ ,  $EU_2(NM|p) = (1-2p)a$ ,  $EU_2(WA|p) = 0$ .

$$EU_{2}(M|p) \ge EU_{2}(NM|p) \quad \Leftrightarrow \quad a^{2} \ge (1-2p)a \qquad \qquad \Leftrightarrow \quad p \ge \frac{1-a}{2}$$
$$EU_{2}(M|p) \ge EU_{2}(WA|p) \quad \Leftrightarrow \quad (1-2p)a + a^{2} \ge 0 \qquad \qquad \Leftrightarrow \quad p \le \frac{1+a}{2}$$

From the two expressions, and given that  $0 < \frac{1-a}{2} < \frac{1+a}{2} < 1$  for  $a \in (0, 1/2)$ , we obtain the best-response of S2 when she believes that S1 commits fraud with probability p to be as described in Expression (5) with  $p_1^* = \frac{1-a}{2}$  and  $p_2^* = \frac{1+a}{2}$ .

Let us now turn our attention to the case  $a \neq b$ . Let us first compare  $EU_2(NM|p)$  with  $EU_2(M|p)$ .

$$EU_2(M|p) \ge EU_2(NM|p) \quad \Leftrightarrow \quad [pa+(1-p)b]^2 \ge -pa+(1-p)b \quad \Leftrightarrow$$
$$\Leftrightarrow \quad p^2(a-b)^2 + p(2ab-2b^2+a+b) + (b^2-b) \ge 0$$

The inequality is always true when -pa + (1-p)b is negative (i.e.  $p > \frac{b}{a+b}$ ), thus it is enough to focus on what happens in the remaining cases (for  $p \le \frac{b}{a+b}$ ).

Solving the quadratic equation  $p^2(a-b)^2 + p(2ab-2b^2+a+b) + (b^2-b) = 0$  with respect to p gives us two real roots:  $\frac{2b(b-a)-(a+b)+\sqrt{(a+b)^2-8ab(b-a)}}{2(b-a)^2}$  and  $\frac{2b(b-a)-(a+b)-\sqrt{(a+b)^2-8ab(b-a)}}{2(b-a)^2}$ . It is straightforward to verify that the discriminant  $(a+b)^2 - 8ab(b-a)$  is strictly positive for all values of  $(a,b) \in (0,1/2)^2$ . Moreover,  $\frac{2b(b-a)-(a+b)-\sqrt{(a+b)^2-8ab(b-a)}}{2(b-a)^2}$  is strictly negative for all values of  $(a,b) \in (0,1/2)^2$ . For 2b(b-a) > (a+b) one should observe that

$$2b(b-a) - (a+b) - \sqrt{(a+b)^2 - 8ab(b-a)} < 0 \Leftrightarrow (b-a)^2(b-1)b < 0$$

which is true for all  $(a, b) \in (0, 1/2)^2$  such that  $a \neq b$ , whereas for  $2b(b-a) - (a+b) \leq 0$  the result is readily obvious. Similarly, it readily follows that  $\frac{2b(b-a)-(a+b)+\sqrt{(a+b)^2-8ab(b-a)}}{2(b-a)^2}$  is strictly positive for all  $(a, b) \in (0, 1/2)^2$ . Finally, turning our attention back to the inequality, given that  $(a-b)^2 > 0$ , the polynomial is positive outside of the roots and negative inside. Thus, overall, we obtain the following result

$$EU_2(M|p) \ge EU_2(NM|p) \quad \Leftrightarrow \quad p \ge p_1^* := \frac{2b(b-a) - (a+b) + \sqrt{(a+b)^2 - 8ab(b-a)}}{2(b-a)^2}$$

with equality holding only at  $p = p_1^*$ .

Let us now compare  $EU_2(M|p)$  with  $EU_2(WA|p)$ .

$$EU_2(M|p) \ge EU_2(WA|p) \quad \Leftrightarrow \quad [pa+(1-p)b]^2 \ge pa-(1-p)b \quad \Leftrightarrow$$
$$\Leftrightarrow \quad p^2(a-b)^2 + p(2ab-2b^2-a-b) + (b^2+b) \ge 0$$

Solving the quadratic equation  $p^2(a-b)^2 + p(2ab-2b^2-a-b) + (b^2+b) = 0$  with respect to p gives two real roots:  $\frac{2b(b-a)+(a+b)-\sqrt{(a+b)^2+8ab(b-a)}}{2(b-a)^2}$  and  $\frac{2b(b-a)+(a+b)+\sqrt{(a+b)^2+8ab(b-a)}}{2(b-a)^2}$ , where it is straightforward to verify that  $(a+b)^2 + 8ab(b-a)$  is always strictly positive in  $(a,b) \in (0,1/2)^2$ . We can also show that the second root is larger than one.

$$2b(b-a) + (a+b) + \sqrt{(a+b)^2 + 8ab(b-a)} > 2(b-a)^2 \iff$$

$$\sqrt{(a+b)^2 + 8ab(b-a)} > 2a(a-b) - (a+b) \qquad \Leftrightarrow \qquad$$

[if rhs positive] 
$$(a+b)^2 + 8ab(b-a) > (2a(a-b) - (a+b))^2 \Leftrightarrow$$
$$2ab(b-a) > a^2(a-b)^2 - a(a^2 - b^2) \Leftrightarrow$$
$$1 > a$$

If the rhs of the third inequality is negative, then the result is trivially true. Similarly, it follows that  $\frac{2b(b-a)+(a+b)-\sqrt{(a+b)^2+8ab(b-a)}}{2(b-a)^2}$ is smaller than one. As far as it concerns the sign of the inequality,
as  $(a-b)^2 > 0$  the polynomial is positive outside the roots and negative inside. Thus, overall

$$EU_2(M|p) > EU_2(WA|p) \quad \Leftrightarrow \quad p \le p_2^* := \frac{2b(b-a) + (a+b) - \sqrt{(a+b)^2 + 8ab(b-a)}}{2(b-a)^2}$$

with equality holding only for  $p = p_2^*$ .

The final step is to compare  $p_1^*$  with  $p_2^*$ .

$$\begin{array}{rcl} p_1^* &<& p_2^* & \Leftrightarrow \\ \\ \frac{2b(b-a) - (a+b) + \sqrt{(a+b)^2 - 8ab(b-a)}}{2(b-a)^2} &<& \frac{2b(b-a) + (a+b) - \sqrt{(a+b)^2 + 8ab(b-a)}}{2(b-a)^2} &\Leftrightarrow \\ \\ \sqrt{(a+b)^2 - 8ab(b-a)} + \sqrt{(a+b)^2 + 8ab(b-a)} &<& 2(a+b) &\Leftrightarrow \\ \\ \sqrt{1 - \frac{8ab(b-a)}{a+b}} + \sqrt{1 + \frac{8ab(b-a)}{a+b}} &<& 2 \end{array}$$

which holds for all  $(a, b) \in (0, 1/2)^2$  such that  $a \neq b$ . To show that the last expression holds, let us first consider the case b > a. Let  $x = \frac{8ab(b-a)}{a+b}$  and observe that, given  $(a, b) \in (0, 1/2)^2$ , for b > a we have  $x \in (0, 1)$ . Then, it is straightforward to see that the function  $f(x) = \sqrt{1-x} + \sqrt{1+x}$  is strictly decreasing in (0, 1) and f(0) = 2. Therefore, for all  $x \in (0, 1)$ , f(x) < 2. Thus,  $\sqrt{1 - \frac{8ab(b-a)}{a+b}} + \sqrt{1 + \frac{8ab(b-a)}{a+b}} < 2$  as well. The case of a > b is identical, simply observing that the expression can be rewritten as  $\sqrt{1 + \frac{8ab(a-b)}{a+b}} + \sqrt{1 - \frac{8ab(a-b)}{a+b}} < 2$ .

Therefore, considering together the expressions: (a)  $EU_2(M|p) \ge EU_2(NM|p) \Leftrightarrow p \ge p_1^*$ , (b)  $EU_2(M|p) \le EU_2(WA|p) \Leftrightarrow p \le p_2^*$ , and (c)  $p_1^* < p_2^*$ , we obtain the best response of S2 as a function of p to be as described in Expression (5) with  $p_1^* = \frac{2b(b-a)-(a+b)+\sqrt{(a+b)^2-8ab(b-a)}}{2(b-a)^2}$  and  $p_2^* = \frac{2b(b-a)+(a+b)-\sqrt{(a+b)^2+8ab(b-a)}}{2(b-a)^2}$ .

Overall, the best-response of S2 to each p is as follows:

$$(q_1^*, q_2^*) = \begin{cases} (0, 1), & \text{if } p < p_1^* \\ (0, 0), & \text{if } p \in (p_1^*, p_2^*) \\ (1, 0), & \text{if } p > p_2^* \\ (0, q_2^* \in [0, 1]), & \text{if } p = p_1^* \\ (q_1^* \in [0, 1], 0), & \text{if } p = p_2^* \end{cases}$$
(5)

where

$$p_1^* = \begin{cases} \frac{2b(b-a) - (a+b) + \sqrt{(a+b)^2 - 8ab(b-a)}}{2(b-a)^2}, & \text{if } a \neq b \\ \frac{1-a}{2}, & \text{if } a = b \end{cases} \qquad p_2^* = \begin{cases} \frac{2b(b-a) + (a+b) - \sqrt{(a+b)^2 + 8ab(b-a)}}{2(b-a)^2}, & \text{if } a \neq b \\ \frac{1+a}{2}, & \text{if } a = b \end{cases}$$
(6)

Best Response of S1: The expected utilities of the S1 for each choice are the following:

$$EU_1(H|q_1, q_2, z) = q_2 R_1 + (1 - q_1 - q_2) \left(\frac{1}{2} + z\right) R_1$$
$$EU_1(F|q_1, q_2, z) = q_2 \underline{R}_1 + (1 - q_1 - q_2) \left(\frac{1}{2} - z\right) \underline{R}_1$$

From the above, it is straightforward to observe that the best-response of **S1** for each  $(q_1, q_2)$  and z is as follows:

$$p = \begin{cases} 1 & \text{if } q_2(\underline{R}_1 - R_1) + (1 - q_1 - q_2) \left( \frac{1}{2}(\underline{R}_1 - R_1) - z(\underline{R}_1 + R_1) \right) > 0 \\ 0 & \text{if } q_2(\underline{R}_1 - R_1) + (1 - q_1 - q_2) \left( \frac{1}{2}(\underline{R}_1 - R_1) - z(\underline{R}_1 + R_1) \right) < 0 \\ [0, 1] & \text{if } q_2(\underline{R}_1 - R_1) + (1 - q_1 - q_2) \left( \frac{1}{2}(\underline{R}_1 - R_1) - z(\underline{R}_1 + R_1) \right) = 0 \end{cases}$$
(7)

The expression can be rewritten with respect to  $k = \frac{1}{2} \frac{\underline{R}_1 - R_1}{\underline{R}_1 + R_1}$  as follows:

$$p = \begin{cases} 1 & \text{if } 2q_2k + (1 - q_1 - q_2)(k - z) > 0 \\ 0 & \text{if } 2q_2k + (1 - q_1 - q_2)(k - z) < 0 \\ [0, 1] & \text{if } 2q_2k + (1 - q_1 - q_2)(k - z) = 0 \end{cases}$$

$$\tag{8}$$

**Proof of Proposition 1:** We test all possible equilibria and find which ones can be sustained and for which parameter values:

- 1.  $(p^* = 0)$  If  $p^* = 0 \Rightarrow q_2^* = 1$  and  $z^* = b$ . If  $q_2^* = 1$  then  $p^* = 1$ . This leads to a contradiction. No equilibrium can be sustained with  $p^* = 0$ .
- 2.  $(p^* = 1)$  If  $p^* = 1$  then  $q_1^* = 1$  and  $z^* = a$ . If  $q_1^* = 1$  then  $p^* \in [0, 1]$ . Therefore, for  $p^* = 1$ , the following equilibrium can be sustained:

$$p^* = 1, (q_1^*, q_2^*) = (1, 0), z^* = a$$

3.  $(p^* \in (0,1) \text{ and } q_1^* = 1:)$  If  $q_1^* = 1$  then  $p^* \in [0,1]$ , but for  $q_1^* = 1$  to be a best response for **S2** it must hold that  $p^* \in [p_2^*, 1]$ . Moreover, for **S1** to be fully mixing it must hold that  $2q_2^*k + (1 - q_1^* - q_2^*)(k - z^*) = 0$ , which holds for  $(q_1^*, q_2^*) = (1, 0)$ . Therefore, the following family of equilibria can be sustained:

$$p^* \in [p_2^*, 1), (q_1^*, q_2^*) = (1, 0), z^* = p^*a + (1 - p^*)b$$

The equilibria of cases 2 and 3 combined give part 1 of the proposition.

- 4.  $(p^* \in (0, 1) \text{ and } q_2^* = 1)$  If  $q_2^* = 1$  then  $p^* = 1$ , which leads to a contradiction.
- 5.  $(p^* \in (0,1) \text{ and } q_1^* = q_2^* = 0)$  If  $(q_1^*, q_2^*) = (0,0)$  then  $p^* \in (0,1)$  can be a best response only if  $z^* = k$ . Thus,  $p^*a + (1 - p^*)b = k$ , which holds if and only if a = b = k or  $a \neq b$ and  $p^* = \frac{k-b}{a-b}$  (if its value is admissible). Moreover,  $(q_1^*, q_2^*) = (0,0)$  is a best response when  $p^* \in [p_1^*, p_2^*]$ . Thus, the following equilibrium can be sustained

$$p^* = \frac{k-b}{a-b}$$
  $(q_1^*, q_2^*) = (0, 0)$  and  $z^* = k$ 

as long as a = b = k, or  $a \neq b$  and  $\frac{k-b}{a-b} \in [p_1^*, p_2^*]$ .

6.  $(p^* \in (0, 1), q_1^* \in (0, 1), \text{ and } q_2^* = 0:)$  For  $q_1^* \in (0, 1)$  and  $q_2^* = 0$  to be (jointly) a best response it must be the case that  $p^* = p_2^*$ . For  $p^* = p_2^*$  to be the best response when  $q_1^* \in (0, 1)$  and  $q_2^* = 0$  it must also be true that  $z^* = k$ . But then this means that  $k = p_2^*a + (1 - p_2^*)b$ , which holds only when a = b = k or  $a \neq b$  and  $p_2^* = \frac{k-b}{a-b}$ . Therefore, the following equilibrium can be sustained:

$$p^* = p_2^*,$$
  $(q_1^*, q_2^*) = (\hat{q}, 0)$  for  $\hat{q} \in (0, 1)$ , and  $z^* = k$ 

as long as a = b = k, or as long as  $a \neq b$  and  $\frac{k-b}{a-b} = p_2^*$ .

7.  $(p^* \in (0,1), q_1^* = 0, \text{ and } q_2^* \in (0,1):)$  For  $q_1^* = 0$  and  $q_2^* \in (0,1)$  to be (jointly) a best response it must be the case that  $p = p_1^*$ . For this to be a best response it would then be required that  $2q_2^*k + (1-q_2^*)(k-z^*) = 0$ , or equivalently  $q_2^* = \frac{z^*-k}{z^*+k} = \frac{p_1^*a + (1-p_1^*)b-k}{p_1^*a + (1-p_1^*)b+k}$  as long as  $k \leq p_1^*a + (1-p_1^*)b$ . In this case, the following equilibrium can be sustained:

$$p^* = p_1^*$$
  $(q_1^*, q_2^*) = \left(0, \frac{p_1^*a + (1 - p_1^*)b - k}{p_1^*a + (1 - p_1^*)b + k}\right)$  and  $z^* = p_1^*a + (1 - p_1^*)b$ 

as long as  $k \le p_1^* a + (1 - p_1^*) b$ 

**Proof of Corollary 1:** For **S2** the result is readily observable from the fact that  $EU_2(M|p)$  is strictly decreasing in p – since  $\frac{\partial EU_2(M|p)}{\partial p} = -a - b + 2(a - b)(pa + (1 - p)b)$ , where  $pa + (1 - p)b \in$ (0, 1/2). Given this, observe that the expected payoff for **S2** in equilibrium 4 (mixed, between Monitor and Not Monitor) is equal to  $EU_2(M|p_1^*)$ , whereas for equilibrium 2 (pure, always Monitor) it is equal to  $EU_2(M|p^*)$  for some  $p^* \in (p_1^*, p_2^*)$ , thus smaller than  $EU_2(M|p_1^*)$ , and for equilibrium 1 (pure, always Walk Away) it is equal to 0, which is by construction also equal to  $EU_2(M|p_2^*)$ , thus again smaller than  $EU_2(M|p_1^*)$ .

For **S1**, in both equilibria 2 and 4 player **S1** uses a mixed strategy, thus both pure strategies yield him the same expected payoff, so it is sufficient to compare his payoffs when honest. In both equilibria, these payoffs are trivially higher than 0 – which is the payoff in equilibrium 1 – given that  $R_1 > 0$ . Between the other two, we need to show that  $(\frac{1}{2} + k) R_1 < \frac{z_1^* - k}{z_1^* + k} R_1 + \frac{2k}{z_1^* + k} (\frac{1}{2} + z_1^*) R_1$ , where  $z_1^* = p_1^* a + (1 - p_1^*) b$ .

$$\left(\frac{1}{2}+k\right)R_1 < \frac{z_1^*-k}{z_1^*+k}R_1 + \frac{2k}{z_1^*+k}\left(\frac{1}{2}+z_1^*\right)R_1 \Leftrightarrow \left(\frac{1}{2}+k\right)(z_1^*+k) < z_1^*+2kz_1^* \Leftrightarrow \left(\frac{1}{2}+k\right)(z_1^*-k) > 0$$

Hence, it is enough that  $z_1^* > k$  or equivalently  $p_1^*a + (1 - p_1^*)b > k$  which is true whenever a > band  $p_1^* > \frac{k-b}{a-b}$ , i.e. in the parameter range in which both equilibrium 1 and equilibrium 4 exist.

**Proof of Proposition 2:** Let us first present the results regarding the effect of a change in k. The result is immediately obtained by differentiating the respective quantities, recalling that  $z^*$  and  $p_1^*$  are independent of k. Thus, first,  $\frac{\partial TR}{\partial k} = z^* \frac{2z^*(1-2p_1^*)-1}{(z^*+k)^2} < 0$  because  $z^* < 1/2$ ,  $1 - 2p_1^* < 1$ . Second,  $\frac{\partial HR}{\partial k} = (1 - p_1^*)z^* \frac{2z^*-1}{(z^*+k)^2} < 0$ , again because  $z^* < 1/2$ . Third,  $\frac{\partial FR}{\partial k} = -p_1^* z^* \frac{1+2z^*}{(z^*+k)^2} < 0$ . Finally,  $\frac{HR}{TR} = 1 - p_1^* \frac{1-2k}{1+2k(1-2p_1^*)}$ . Thus,  $\frac{\partial (HR/TR)}{\partial k} = 4p_1^* \frac{1-p_1^*}{[1+2k(1-2p_1^*)]^2} > 0$ .

Subsequently, we are interested in the derivatives of the relevant quantities with respect to a and b. To obtain the sign of these derivatives, we need some intermediate results, which we prove below.

#### Lemma 1

- $\frac{\partial p_1^*}{\partial a} < 0$ ,  $\frac{\partial z^*}{\partial a} > 0$ , and  $\frac{\partial q_2^*}{\partial a} > 0$
- $\frac{\partial p_1^*}{\partial b} > 0$ ,  $\frac{\partial z^*}{\partial b} > 0$ , and  $\frac{\partial q_2^*}{\partial b} > 0$

**Proof of Lemma 1:** Let us first prove that  $\frac{\partial p_1^*}{\partial a} < 0$ . It is helpful to consider the implicit form of obtaining  $p_1^*$ , i.e. the equation  $EU_2(M|p_1^*) = EU_2(NM|p_1^*)$  or equivalently  $p_1^*a - (1 - p_1^*)b + [p_1^*a + (1 - p_1^*)b]^2 = 0$ . If we differentiate this expression implicitly with respect to a we

get  $\frac{\partial p_1^*}{\partial a} \{a+b+2(a-b) [p_1^*a+(1-p_1^*)b]\} + p_1^* [1+2p_1^*a+2(1-p_1^*)b] = 0$ . Thus, given that both  $1+2p_1^*a+2(1-p_1^*)b>0$  and  $a+b+2(a-b) [p_1^*a+(1-p_1^*)b] > 0$ ,<sup>12</sup> it follows that  $\frac{\partial p_1^*}{\partial a} < 0$ .

Similarly, we shall use the same expression to show that  $\frac{\partial p_1^*}{\partial b} > 0$ . Namely, by implicitly differentiating  $p_1^*a - (1 - p_1^*)b + [p_1^*a + (1 - p_1^*)b]^2 = 0$  with respect to b we get the following expression:  $\frac{\partial p_1^*}{\partial b} \{a + b + 2(a - b) [p_1^*a + (1 - p_1^*)b]\} + (1 - p_1^*) [2p_1^*a + 2(1 - p_1^*)b - 1] = 0$ . From this, note that  $2p_1^*a + 2(1 - p_1^*)b - 1 < 0$  for all  $(a, b) \in (0, 1/2)^2$ , which combined with the fact that  $a + b + 2(a - b) [p_1^*a + (1 - p_1^*)b] > 0$ , which was proven before, implies that  $\frac{\partial p_1^*}{\partial b} > 0$ .

Let us now prove that  $\frac{\partial z^*}{\partial a} > 0$ . To do so, it is useful to observe that if we combine  $p_1^*a - (1-p_1^*)b + [p_1^*a + (1-p_1^*)b]^2 = 0$  and  $z^* = p_1^*a + (1-p_1^*)b$ , we can obtain the equation:  $z^* + (z^*)^2 - 2(1-p_1^*)b = 0$ . If we differentiate this expression implicitly with respect to a, and after some rearrangements, we get that  $\frac{\partial z^*}{\partial a} = -\frac{2b}{1+2z^*}\frac{\partial p_1^*}{\partial a} > 0$ , given that we have already shown that  $\frac{\partial p_1^*}{\partial a} < 0$ .

Similarly to this, expression  $z^* + (z^*)^2 - 2(1-p_1^*)b = 0$  can be rewritten as  $2p_1^*a - z^* + (z^*)^2 = 0$ . Differentiating this implicitly with respect to b we get that  $\frac{\partial z^*}{\partial b} = \frac{2a}{1-2z^*}\frac{\partial p_1^*}{\partial b} > 0$ , given that we have already shown that  $\frac{\partial p_1^*}{\partial b} > 0$ .

Finally,  $\frac{\partial q_2^*}{\partial a} = \frac{\partial \left(\frac{z^*-k}{z^*+k}\right)}{\partial a} = \frac{2k}{(z^*+k)^2} \frac{\partial z^*}{\partial a} > 0$  and  $\frac{\partial q_2^*}{\partial b} = \frac{\partial \left(\frac{z^*-k}{z^*+k}\right)}{\partial b} = \frac{2k}{(z^*+k)^2} \frac{\partial z^*}{\partial b} > 0$ , given that  $\frac{\partial z^*}{\partial a} > 0$  and  $\frac{\partial z^*}{\partial b} > 0$  respectively.

Let us now look at the derivatives with respect to a. First, regarding total research, we observe that  $\frac{\partial HR}{\partial a} = (1+2k) \left[ \frac{k}{(z^*+k)^2} (1-p_1^*) \frac{\partial z^*}{\partial a} - \frac{z^*}{z^*+k} \frac{\partial p_1^*}{\partial a} \right] > 0$  given that  $\frac{\partial p_1^*}{\partial a} < 0$  and  $\frac{\partial z^*}{\partial a} > 0$ . Second, regarding fraudulent research, recalling that  $\frac{\partial z^*}{\partial a} = -\frac{2b}{1+2z^*} \frac{\partial p_1^*}{\partial a}$  we observe that  $\frac{\partial FR}{\partial a} = (1-2k) \left[ \frac{k}{(z^*+k)^2} p_1^* \frac{\partial z^*}{\partial a} + \frac{z^*}{z^*+k} \frac{\partial p_1^*}{\partial a} \right] = \frac{(1-2k)}{(z^*+k)^2(1+2z^*)} \frac{\partial p_1^*}{\partial a} [z^*(z^*+k)(1+2z^*) - 2bkp_1^*]$ . Let us now focus on the term  $z^*(z^*+k)(1+2z^*) - 2bkp_1^*$ . If we use the expression  $z^* + (z^*)^2 - 2(1-p_1^*)b = 0$  that we obtained above, we can rewrite this expressions as  $2(z^*)^3 + (3k+1)(z^*)^2 + 2kz^* - 2bk$ . Knowing that  $\frac{\partial z^*}{\partial a} > 0$ , this expression increases in a and thus changes sign at most once as a increases. Let us first see what happens at a = 0. Therefore, for a = 0,  $p_1^* = 1$  and  $z^* = 0$ , the expression becomes equal to -2bk < 0. On the other hand, for a = b,  $p_1^* = \frac{1-b}{2}$  and  $z^* = b$ , so the expression becomes equal to  $2b^3 + (3k+1)b^2 > 0$ . Hence, there exists some  $\hat{a}(b) \in (0, b)$  such that  $\frac{\partial FR}{\partial a} > 0$  for  $a > \hat{a}(b)$ .

<sup>&</sup>lt;sup>12</sup>To see that  $a + b + 2(a - b) [p_1^*a + (1 - p_1^*)b] > 0$  observe that the expression increases in  $p_1^*$ , thus it is enough to show that it is positive for  $p_1^* = 0$ , where the expression becomes a + b + 2(a - b)b. This increases in a, thus it is enough to check for a = 0, where this becomes  $b - 2b^2$ , which is strictly larger than 0 for all  $b \in (0, 1/2)$ .

Yet, it still remains unclear whether this threshold falls within the parameter region in which the mixed equilibrium exists, i.e. where  $z^* \ge k$ . Therefore, we need to check for which values of k, the inequalities  $z^* \ge k$  and  $z^*(z^* + k)(1 + 2z^*) - 2bkp_1^* \le 0$  hold simultaneously for some values of a and b. This is the case when  $k \in \left[\frac{(z^*)^2(1+2z^*)}{2bp_1^*-z^*(1+2z^*)}, z^*\right]$ , as long as  $z^* \ge \frac{(z^*)^2(1+2z^*)}{2bp_1^*-z^*(1+2z^*)}$  or equivalently  $z^*(1 + 2z^*) \le bp_1^* \Leftrightarrow a \le \frac{3b\sqrt{9+40b}-b(7+8b)}{4(1+8b)}$ . The maximum  $z^*$  for which the inequality holds is equal to  $\frac{\sqrt{29}-3}{10}$ . Thus, for all  $k > \hat{k} = \frac{\sqrt{29}-3}{10}$  the two inequalities never hold together, which in turn means that the threshold always falls outside the parameter region in which the mixed equilibrium exists. On the contrary, for all b, as a tends to zero, the lower bound  $\frac{(z^*)^2(1+2z^*)}{2bp_1^*-z^*(1+2z^*)}$ tends to zero. Hence, for all  $k \in (0, \hat{k})$ , for all  $b \in (0, 1/2)$ , there exist a for which both inequalities hold simultaneously.

Third,  $\frac{\partial TR}{\partial a} = \frac{\partial HR}{\partial a} + \frac{\partial FR}{\partial a} = \frac{k}{(z^*+k)^2} \left[ (1+2k)(1-p_1^*) + (1-2k)p_1^* \right] \frac{\partial z^*}{\partial a} - 4k \frac{z^*}{z^*+k} \frac{\partial p_1^*}{\partial a} > 0$  given that  $\frac{\partial z^*}{\partial a} > 0$ ,  $\frac{\partial p_1^*}{\partial a} < 0$ , and  $(1+2k)(1-p_1^*) + (1-2k)p_1^* > 0$ . Finally,  $\frac{\partial (HR/TR)}{\partial a} = -\frac{(1-2k)(1+2k)}{\left[(1+2k)-4kp_1^*\right]^2} \frac{\partial p_1^*}{\partial a} > 0$ , because  $\frac{\partial p_1^*}{\partial a} < 0$ .

Looking now at the derivatives with respect to b, we get the following results. Regarding FR,  $\frac{\partial FR}{\partial b} = (1-2k) \left[ \frac{k^*}{(z^*+k)^2} p_1^* \frac{\partial z^*}{\partial b} + \frac{z^*}{z^*+k} \frac{\partial p_1^*}{\partial b} \right] > 0$ , given that  $\frac{\partial z^*}{\partial b} > 0$  and  $\frac{\partial p_1^*}{\partial b} > 0$ . The impact of b on HR and TR is less clear, but not for the ratio HR/TR. Namely, observe that  $\frac{HR}{TR} = \frac{1+2k-(1+2k)p_1^*}{1+2k-4kp_1^*}$ . Thus,  $\frac{\partial (HR/TR)}{\partial b} = -\frac{(1+2k)(1-2k)}{(1+2k-4kp_1^*)^2} \frac{\partial p_1^*}{\partial b} < 0$ , given that  $\frac{\partial p_1^*}{\partial b} > 0$ .

**Proof of Proposition 3:** First, we show that HR - FR increases in *a* for any *k* and *b*. It holds, because  $HR - FR = TR\left(2\frac{HR}{TR} - 1\right)$  and recall from Proposition 2 that TR and HR/TR increase in *a*, while FR/TR decreases in *a*. Therefore, the objective is maximized at the maximum value of *a*.

For the optimal k, let us rewrite the objective function as follows:  $HR - FR = \frac{z^*}{z^*+k} (1 + 2k - 2p_1^*)$ . Analyzing the monotonicity of this quantity with respect to k, we observe that the optimal k depends on the values of a and b, but in any case it is either equal to  $k_{min}$  or to  $k_{max}$ . More specifically,  $\frac{\partial(HR-FR)}{\partial k} = \frac{z^*}{(z^*+k)^2}(2z^* - 1 + 2p_1^*)$ , which for any  $a \in (0, 1/2)$  is negative for b = 0 and positive for b = 1/2. Thus, for each (a, b) it is either positive for any k or negative for any k, and there is a threshold  $\hat{b}$  which separates the two domains.

Therefore, we obtain that for low values of b, the function is maximized at  $(a, k) = (a_{max}, k_{min})$ , whereas for high values of b, it is maximized at  $(a, k) = (a_{max}, k_{max})$ .

# References

- DE QUIDT, J., T. FETZER, AND M. GHATAK (2016): "Group lending without joint liability," Journal of Development Economics, 121, 217–236.
- FANELLI, D., R. COSTAS, AND V. LARIVIÈRE (2015): "Misconduct policies, academic culture and career stage, not gender or pressures to publish, affect scientific integrity," *PloS one*, 10, e0127556.
- GALL, T., J. P. IOANNIDIS, AND Z. MANIADIS (2017): "The credibility crisis in research: Can economics tools help?" *PLoS biology*, 15, e2001846.
- GHATAK, M. AND T. W. GUINNANE (1999): "The economics of lending with joint liability: theory and practice," *Journal of development economics*, 60, 195–228.
- HALL, K. L., A. L. VOGEL, G. C. HUANG, K. J. SERRANO, E. L. RICE, S. P. TSAKRAKLIDES, AND S. M. FIORE (2018): "The science of team science: A review of the empirical evidence and research gaps on collaboration in science." *American psychologist*, 73, 532.
- HELGESSON, G. AND S. ERIKSSON (2018): "Responsibility for scientific misconduct in collaborative papers," *Medicine, Health Care and Philosophy*, 21, 423–430.
- HUSSINGER, K. AND M. PELLENS (2019): "Scientific misconduct and accountability in teams," *Plos one*, 14, e0215962.
- KARLAN, D. S. (2007): "Social connections and group banking," The Economic Journal, 117, F52–F84.
- KIRI, B., N. LACETERA, AND L. ZIRULIA (2018): "Above a swamp: A theory of high-quality scientific production," *Research Policy*, 47, 827–839.
- LACETERA, N. AND L. ZIRULIA (2011): "The economics of scientific misconduct," *The Journal* of Law, Economics, & Organization, 27, 568–603.
- LARIVIÈRE, V., Y. GINGRAS, C. R. SUGIMOTO, AND A. TSOU (2015): "Team size matters: Collaboration and scientific impact since 1900," Journal of the Association for Information Science and Technology, 66, 1323–1332.

- MARX, L. M. AND F. SQUINTANI (2009): "Individual accountability in teams," Journal of Economic Behavior & Organization, 72, 260–273.
- MICELI, T. J. AND K. SEGERSON (2007): "Punishing the innocent along with the guilty: The economics of individual versus group punishment," *The Journal of Legal Studies*, 36, 81–106.
- POLDRACK, R. A., C. I. BAKER, J. DURNEZ, K. J. GORGOLEWSKI, P. M. MATTHEWS, M. R. MUNAFÒ, T. E. NICHOLS, J.-B. POLINE, E. VUL, AND T. YARKONI (2017): "Scanning the horizon: towards transparent and reproducible neuroimaging research," *Nature reviews neuroscience*, 18, 115–126.
- THORP, H. H. (2023): "Generative approach to research integrity," .
- VIE, K. J. (2020): "How should researchers cope with the ethical demands of discovering research misconduct? Going beyond reporting and whistleblowing," *Life Sciences, Society and Policy*, 16, 1–18.