

# Vague by Design: Performance Evaluation and Learning from Wages\*

Franz Ostrizek<sup>†</sup>

December 4, 2019

## Abstract

We study a dynamic principal-agent setting in which both sides learn about the importance of effort. The quality of the agent’s output is not observed directly. Instead, the principal jointly designs an evaluation technology and a wage schedule. More precise performance evaluation reduces current agency costs but promotes learning, which is shown to increase future agency costs. As a result, the optimal evaluation technology is both imprecise and tough: a bad performance is always sanctioned, but a good one is not always recognized.

We also study the case in which principal and agent have different priors, for instance because the agent has incorrect beliefs about his abilities. If the agent is overconfident, the principal uses a tough evaluation structure to preserve the agent’s profitable misperception. For an underconfident agent, by contrast, she either uses a fully informative evaluation in order to promote learning and eliminate costly underconfidence, or is lenient if learning is too costly.

## 1 Introduction

Many firms motivate their workers to exert effort with incentive pay based on objective measures of performance.<sup>1</sup> Such measures have become richer and easier

---

\*I am grateful to Roland Bénabou, Pietro Ortoleva, Wolfgang Pesendorfer for their continued guidance and to Ludmila Matysková, Sofia Moroni, Leon Musolff, Pellumb Reshidi, Evgenii Safonov, Denis Shishkin, Nikhil Vellodi, Can Urgan, Leeat Yariv and seminar audiences at Princeton University for helpful comments and discussions.

<sup>†</sup>Princeton University. ostrizek@princeton.edu

<sup>1</sup>According to data from the Bureau of Labor Statistics’ National Compensation Survey, 39% of hours worked in US private sector firms in 2013 were in jobs with performance-related pay. 21%

to obtain as a result of the IT revolution. For example, the availability of board computers and GPS tracking allows for better monitoring of truck drivers and time tracking software in law firms and other offices not only simplifies billing but also logs the activities of employees. Shop floor control systems monitor not only the flow of goods but also allow the tracking of workers. Improved natural language processing enables data collection in applications ranging from call centers to health care.<sup>2</sup> Based on such statistics, firms can arrive at better objective measures of workers' contribution to profits.

Should this additional information be used to set performance pay? The theory of incentives seems to offer a simple answer to this question. Providing incentives to workers is costly and may cause inefficiencies because the underlying performance measures are only partially informative about effort. Conversely, more hard information about the workers' contribution is always helpful and should be used as a basis of performance pay (Holmström, 1979; Grossman and Hart, 1983). In particular, it is never optimal to base incentives on a noisy signal instead of the contribution to output.<sup>3</sup>

In this paper, we show that learning changes this conclusion fundamentally. The firm prefers workers to remain uninformed about their match-specific ability since uninformed workers are cheaper to motivate on average. Contingent wages, however, reveal information about output and hence their ability. As a result, the principal wants to base incentive pay not on the agent's performance directly but on a noisy signal instead.

To illustrate how explicit incentives invariably reveal information, suppose a worker's contract promises a bonus upon a sufficiently high customer satisfaction score. The worker does not observe this score directly. If he receives the bonus, however, he knows that he cleared the threshold while if he does not receive it, he learns that he fell short of it.

---

fall into a narrower classification of performance-related pay excluding, among other categories, referral bonuses which should arguably be excluded from a theoretical perspective, but also safety bonuses which should be retained (Gittleman and Pierce, 2013).

<sup>2</sup>These tools allow call centers to detect the callers' mood, for example (Singer, 2013). For a survey on the use of natural language processing to extract information from health-related text, see Gonzalez-Hernandez et al. (2017).

<sup>3</sup>That is not optimal to add noise continues to hold across a wider class of models, including multi-tasking (Holmström and Milgrom, 1991) and linear-normal career concerns models Holmström (1999); Hörner and Lambert (forthcoming).

Why would it be profitable to conceal information about their performance from workers? The costs of providing such incentives itself provides a rationale. When ability and effort are complements in production, an agent who believes his ability is high only requires a small bonus to motivate him to exert high effort and vice versa. The impact of a given change in beliefs is amplified if it is large relative to the expected efficiency of effort. Therefore, it has a large impact at a low posterior, while it has a smaller impact at a high posterior. This implies that learning, which causes a mean-preserving spread of posterior beliefs, is costly on average: at low beliefs, the required bonus increases a lot, while at high beliefs, the required bonus decreases only a little.<sup>4</sup>

To capture this fundamental trade-off between incentives and information, we develop a model of twice-repeated moral hazard with learning. The agent’s type affects not only mean output, but also the effectiveness of effort. The quality of output is not observed directly. Instead, in every period the principal designs not only wages but also the underlying performance evaluation. The agent observes his evaluations and wages. The evaluation structure therefore determines not only the cost of incentives this period, but also the extent of learning.<sup>5</sup> We analyze the model in Section 3, transforming the contracting problem into an information design problem with additional constraints (participation and incentive compatibility) and an additional choice variable, the wage at every posterior.

In the final period, when the continuation belief is of no importance to the principal, the optimal evaluation structure is fully informative. In the first period, however, there is a novel tradeoff. A more precise evaluation structure reduces agency costs this period, but induces more learning, thereby increasing agency costs in the next period. We solve for the optimal contract and show that it features a *binary evaluation structure*: The additional motive of shaping learning does not add complication relative to the fully informative evaluation. The optimal evaluation structure is “*tough*”: The agent obtains a high evaluation and therefore a bonus only if his output is high quality. Even if output was high, however,

---

<sup>4</sup>Since ability also affects the baseline probability of high output, the exact condition is slightly more complicated and is implied by log-supermodularity, as will be discussed later.

<sup>5</sup>Of course, this mechanism is predicated on two background conditions: First, the agent observes his wages. In particular, it is not possible to record the agents performance measures and provide incentives only at the end of the employment relationship. Second, information that is not used as the basis of explicit incentives does not have to be revealed to the agent in any other way. We discuss these issue in more detail after introducing the model formally in Section 2 and provide extensions in Section 5.

he may receive a low evaluation and thus fail to obtain the bonus. After low quality output, the latter outcome is guaranteed. This information structure avoids inducing very low posteriors. Agents with such beliefs would be very expensive to motivate in the next period and even a small increase in their posterior belief has a large (decreasing) effect on the required bonus.

The influence of explicit incentives on the agent’s learning becomes even more essential when the agent is initially misguided. Indeed, the learning environment then also shapes the evolution of mean beliefs, making it important for the principal to preserve profitable worker misconceptions and eliminate costly ones, and for the analyst to determine the persistence of such misconceptions. In particular, a substantial empirical and experimental literature suggests that people are often overconfident about their ability, the degree of control they have over their environment, or the extent to which they live in a “just world” that rewards effort in the long run. In Section 4, we therefore analyze the model with heterogeneous beliefs, allowing the agent to be optimistic or pessimistic about his type.<sup>6</sup> We show that a noisy and tough information structure remains optimal in the face of overconfidence: The best way to preserve profitable optimism is not “coddling” grade inflation, but tough evaluation. If the agent is pessimistic, the principal is still averse to a dispersion in beliefs, but wants to eliminate costly pessimism: If the latter effect dominates, the principal uses a fully informative evaluation to promote learning. If the latter effect dominates, the optimal evaluation structure is now *lenient*: Sometimes a bad performance nonetheless receives a good evaluation and is rewarded with a bonus.

In Section 5 we consider several extensions of our model. We show that the optimal evaluation remains partially informative and tough if the principal can acquire private information about the agent’s performance, if effort is unobserved in addition to being noncontractible, and when the principal can commit across time periods. Section 6 concludes. The proofs not given in the text are collected in the Appendix.

---

<sup>6</sup>We retain the assumption that the agent is Bayesian. There is some evidence suggesting that individuals update their initially optimistic beliefs about their self-control in such a fashion [Yaouanq and Schwardmann \(2019\)](#).

## Related Literature

This paper contributes to the large literature on information in moral-hazard models, focusing on firm-worker relationships. We offer a counterpoint to the classic results establishing that more precise information reduces agency costs (Holmström, 1979; Grossman and Hart, 1983; Kim, 1995) by providing a setting in which the principal prefers to base wages on a noisy information structure.

Coarse or noisy evaluations are also found by some other literatures. Several contributions show that obtaining the right measure of performance is difficult and using the wrong measure can backfire. This is the case with multitasking (Holmström and Milgrom, 1991) or when the agent has private information that would allow him to game a deterministic incentive scheme (Ederer et al., 2018). Hence, it can be optimal to leave information on the table and even introduce noise into the contract. Another strand of the literature shows that when evaluation is subjective, i.e. based on unverifiable private information of the principal, the resulting equilibrium remuneration will depend only on coarse information (MacLeod, 2003; Fuchs, 2007). We show that noisy evaluation is optimal even if verifiable information about the agent's true performance would be available.

That more, symmetric information about the technology can reduce profits in a moral hazard setting has been noted in the literature in several settings. It is well understood that ex-post incentive compatibility is more demanding than ex-ante incentive compatibility when implementing a fixed action. Lizzeri et al. (2002) show that interim performance evaluation is not optimal when there is no learning.<sup>7</sup> Nafziger (2009) demonstrates that it can be optimal to conceal information until after the agent's effort choice, even though this precludes the principal from adjusting the implemented action. Indeed, such situations are generic if the problem is sufficiently rich (Jehiel, 2015). In all these papers, the wage is still allowed to depend on the true realization of the signal, even if it is not revealed ex-ante. We show that less information about the technology increases profits even if this implies that the wage cannot depend on the state even ex-post.<sup>8</sup>

---

<sup>7</sup>In a tournament setting with exogenous and relative payments depending on cumulative output, the optimality of interim performance evaluation depends on the shape of the effort cost function (Ederer, 2010).

<sup>8</sup>Under this assumption, Fang and Moscarini (2005) show that information is detrimental if it erodes profitable overconfidence, see below.

To our knowledge, this is the first paper to combine the three key features of explicit incentives, learning about a persistent type, and information design. There are several literatures combining each two of these features:

A growing literature investigates the design of information structures in one-shot moral hazard problems with commitment to a wage scheme. The older literature (Dye, 1986; Feltham and Xie, 1994; Datar et al., 2001) considers the optimal acquisition and aggregation of information within a parametric class.<sup>9</sup> In Georgiadis and Szentes (forthcoming) and Li and Yang (forthcoming) the costs of information acquisition are assumed as part of the technology. Hoffmann et al. (2019) analyze a setting where the agent takes a single action, but information about his performance arrives over time. Information acquisition requires delayed payments, which creates endogenous costs because of impatience and imperfect risk sharing. We consider the design of the optimal information structure in a repeated setting with learning. The monitoring of output within every period affects not only incentives but also continuation beliefs, which leads to endogenous costs of information.

Learning about a persistent state and information design are combined in a growing literature. Most closely related to our moral-hazard setting are Smolin (2017) and Ely and Szydlowski (2019). In these papers, the worker updates his beliefs about the value of continuing the employment relationship compared to quitting based on information designed by the principal. The principal uses only information, which is valuable for the agent, as an incentive: she cannot commit to contingent payments. We analyze the role of information design when the principal also designs wages to provide incentives. Since the principal sets incentives and ability is match-specific, (symmetric) information itself is not valuable. Information and incentive design constrain each other, as the principal reveals at least as much information as is contained in wages.

Information and implicit incentives are also linked in models of career concerns (Holmström, 1999). A key difference is that in our paper, the principal commits to a wage payment based on the monitoring outcome, while in models of career concerns wages are determined by the belief of an observer (“the market”). As a consequence, the role of information is different. For career concerns, it is essential

---

<sup>9</sup>Indeed, when restricting attention to linear contracts, it can be optimal to leave information unused. (Feltham and Xie, 1994; Datar et al., 2001) This is a consequence of the restricted space of contracts, however.

that skill and effort jointly affect the performance - the agent is motivated to exert effort because a decrease in output would be interpreted as low skill by the potential employers. Hörner and Lambert (forthcoming) analyze the optimal design of the information structure in a Gaussian career concerns model and show how it combines information from different sources or vintages to achieve the optimal combination of dependence on effort and on the agent’s type.<sup>10</sup> In our setting with explicit incentives, entangled information is the source of the friction: the principal would prefer to provide incentives based on a signal that is informative about effort but independent of the agent’s type in order to prevent learning.

The literature on learning in moral hazard models (Adrian and Westerfield, 2009; Giat et al., 2010; Prat and Jovanovic, 2014; Demarzo and Sannikov, 2017) studies learning based on output instead of an information structure that is designed endogenously by the principal. A key concern in this literature is belief manipulation – the fact that the agent’s belief is private information after he deviated from the proposed effort. We analyze this issue as an extension and show that the basic shape of the optimal information structure is preserved.

Another important distinction is that we consider learning about the importance of effort, as opposed to learning about a state that affects only the level of output. The latter is often considerably more tractable and several of our assumptions are a result of handling this feature.<sup>11</sup> Notable exceptions are Bhaskar and Mailath (2019), who show that with learning about the importance of effort and spot contracts, the costs of implementing effort diverge as the number of periods grows, by establishing lower bounds on the cost of incentives.

Our extension to heterogeneous beliefs connects to the literature on contracting with overconfident agents, in particular de la Rosa (2011) who shows that overconfidence about the impact of effort relaxes the incentive constraint and is profitable for the principal. Fang and Moscarini (2005) show that if workers are sufficiently overconfident, the principal wants to conceal her private information about their

---

<sup>10</sup>They also show that is never optimal to introduce noise into the evaluation when implementing the highest effort. Dewatripont et al. (1999) consider a one-shot career concerns problem when effort and the agent’s type enter output in a general form and show that it can be optimal to add noise to the signal of his performance, as noise may increase the impact of effort on the realized signal.

<sup>11</sup>Importantly in our case, the information design problem solved by the principal cannot be written as the design of a distribution of posterior means unless the agent’s type is binary.

true type by offering the same wage contract (which involves a fully informative evaluation of their output) to all workers. We derive how the principal shapes the performance evaluation to shape learning and preserve this misperception.

Technically, our paper relates to the literature on information design (Kamenica and Gentzkow, 2011; Bergemann and Morris, 2019), in particular the recent development of tools to handle information design problems with constraints (Boleslavsky and Kim, 2017; Le Treust and Tomala, 2019; Doval and Skreta, 2018) and additional choice variables (Georgiadis and Szentes, forthcoming) – in our case wages. We apply these tools in a setting where the information designer chooses a signal structure about one variable – output – in order to affect beliefs about another – the ability of the agent. This feature is particularly important in our extension to heterogeneous beliefs: Even though the prior of the principal and the agent are not mutually absolutely continuous, the information design problem is analyzed as if they were, using the transformation approach of Alonso and Câmara (2016).

## 2 The Model

A principal (she) employs an agent (he) for two periods. The principal is risk neutral, the agent is risk averse with utility index  $u : [0, \infty) \rightarrow [0, \infty)$  which we assume to be unbounded.<sup>12</sup> Both share a common discount factor  $\delta \in (0, 1]$ .

**Technology** Each period, the agent exerts nonverifiable effort  $e_t \in \{0, 1\}$  at cost  $ce_t$ , with  $c > 0$ . The worker has a time-invariant ability  $\theta \in \{\theta_L, \theta_H\}$ . For the main sections, we assume that the principal and the agent share a common prior belief  $\mu$  that the agent has a high ability.

The resulting output has either high or low quality,  $y \in \{y_L, y_H\}$ . We normalize the expected revenue from low output to zero and denote the expected revenue from high output by  $Y > 0$ . The probability of a high quality depends on the agent's effort and type, as follows:

---

<sup>12</sup>This is for simplicity to avoid corner solutions.



\ effort	$e_t = 0$	$e_t = 1$
type		
$\theta = \theta_L$	$a$	$a + b$
$\theta = \theta_H$	$a + \Delta a$	$a + b + \Delta a + \Delta b$

Effort and ability are both productive,  $b \geq 0$  and  $\Delta a \geq 0$ , and the technology is log-supermodular,  $a\Delta b - b\Delta a > 0$ . We assume that the principal wants to implement high effort in both periods and after all histories.<sup>13</sup>

**Information, Contracts and Commitment** We assume that the principal has full commitment within each period, but no commitment across periods. Within every period, timing is as follows: The principal proposes a contract, comprising a signal space  $S$ , a distribution over signals conditional on output, and a mapping from signals to wages.<sup>14</sup> Having observed the contract, the agent decides whether to quit and obtain outside utility  $U$  or to work, choosing effort level  $e_t$ . The outside utility is independent of the agent's type, which is assumed to be match specific, and satisfies  $U > \frac{a}{b}c$ .<sup>15</sup> At the end of the period, output, signals and wages realize.

Output is informative about the agent's type, but not directly observed by the agent or the principal.<sup>16</sup> The principal and the agent observe (noncontractible) effort, signals and wages and update their beliefs about the agent's type according to Bayes rule. Therefore, the evaluation designed by the principal has the dual role of providing the basis for incentive pay and determining the learning environment.

---

<sup>13</sup>It is easy to see that implementing high effort after all histories is optimal for the principal for a sufficiently high gain from high quality output,  $Y$ . This sharpens the trade-off between incentives and learning we aim to investigate, as the principal derives no instrumental value of information. Furthermore, implementing a given effort level is a standard focus in the contracting literature (e.g. [Dittmann and Maug, 2007](#); [Edmans and Gabaix, 2011](#)).

<sup>14</sup>This restriction to deterministic wages conditional on the signal is without loss, as the principal can simply extend the signal space to generate any desired randomness in the wage.

<sup>15</sup>The condition on the outside utility assures that the non-negativity constraint implicit in the utility function is never binding in the optimal contract.

<sup>16</sup>We can allow the firm to observe aggregate quantities. In a large organization, it is difficult to link a shortfall in aggregate outcomes to the individual worker, however. Formally, consider the model with a continuum of agents. Through its regular accounting activities, the firm observes aggregate outcomes such as profits, revenues or the average quality of output. These outcomes are not informative about the performance about an individual, infinitesimal agent.

**Discussion** As mentioned in the introduction, it is important that the firm cannot engage in complete backloading of information while still providing incentives.<sup>17</sup> If the principal can record the output of the agent without revealing this information and credibly commit to contingent payments at the end of the relationship, a fully informative and fully delayed evaluation would be optimal. Our mechanism comes into effect when such complete backloading is costly or infeasible. In the main sections, the lack of intertemporal commitment ensures that incentives for effort have to be provided in the concurrent period. Wages are thus informative about output. We allow the principal to reveal more information through the evaluation, what is crucial, however, is that the agent observes at least as much information as contained in the wages. We discuss other options that preclude full “informational backloading” in Section 5.3, such as an impatient agent or a chance of information leaks.

A model of learning based on an endogenous signal distribution with hidden actions has the potential to create subtle issues of endogenous private information, both for the principal and the agent. For tractability and to focus on the main trade-off between incentives and learning, our assumptions ensure that no endogenous private information arises.

In the benchmark model, the principal does not acquire private information about the type of the agent, since she does not privately observe the quality of output directly, but only through the evaluation structure.<sup>18</sup> We relax this assumption in Section 5.1 and allow the principal to acquire noncontractible private information about the agent’s performance in addition. We show that is optimal for the principal not to acquire such private information. To use it, it would have to affect the contract and would thereby be revealed to the agent in the second period. Then, however, it could have been revealed as part of the evaluation in the first period, in a way that reduced agency costs. Even if the choice of private information by the principal is unobserved, this outcome is supported as the unique equilibrium in a natural class.

On the agent side, the benchmark model ensures that his posterior belief remains common knowledge after a deviation to lower effort, since effort is

---

<sup>17</sup>This is an assumption about the flow of information, the timing of the actual monetary transfer is not crucial.

<sup>18</sup>We can accommodate the principal observing the firm’s aggregate output. Suppose there is a continuum of workers as described. Then, the aggregate outcome is uninformative about an individual agent’s effort and type.

noncontractible but observed. If effort is unobserved, the agent acquires private information about his belief after a deviation and double deviations to low effort in both periods will be strictly profitable. In Section 5.2, we derive the resulting dynamic incentive compatibility constraint, analyze this extended model, and show that our results generalize to this case.

### 3 Analysis

A contract is a tuple  $(S, p, w)$ , where  $S$  is an arbitrary measurable signal space,  $p(\cdot|y) \in \Delta(S)$ <sup>19</sup> denotes the distribution of the signal conditional on high (resp. low) output and  $w : S \rightarrow \mathbb{R}$  denotes the wage promised to the agent after each signal.

#### 3.1 Transformation to Belief Space

Every signal  $s \in S$  induces a posterior belief

$$\mu(s) = \mu \frac{p(s|y_L) + (a + b + \Delta a + \Delta b) [p(s|y_H) - p(s|y_L)]}{p(s|y_L) + (a + b + (\Delta a + \Delta b)\mu) [p(s|y_H) - p(s|y_L)]} \quad (1)$$

by Bayes rule. Note that (1) relies on the presumption that high effort was exerted and is therefore only valid if there is no deviation from the effort proposed in the contract. The posterior is increasing in the likelihood ratio of the signal,  $\frac{p(s|y_H)}{p(s|y_L)}$ , since the high type is more likely to produce high output, and is fully determined by this likelihood ratio. It is bounded between  $\underline{\mu}$  and  $\bar{\mu}$ , where

$$\underline{\mu} = \mu \frac{1 - (a + b + \Delta a + \Delta b)}{1 - (a + b + (\Delta a + \Delta b)\mu)}; \quad \bar{\mu} = \mu \frac{a + b + \Delta a + \Delta b}{a + b + (\Delta a + \Delta b)\mu} \quad (2)$$

denote the posteriors associated to a signal that realizes only after a low output ( $p(s|y_H) = 0$ ) and only after a high output ( $p(s|y_L) = 0$ ), respectively.

The contract in period  $t$  affects profits in that period but also the distribution over posteriors, which determines the continuation value of the principal. In the last period, this value is of course zero. In the first period, it is given by the

---

<sup>19</sup>Throughout, we will use integral notation for expected values and understand expressions of the form  $\int f(x) dx$  in the sense of distributions where required; no absolute continuity with respect to Lebesgue measure is assumed. With slight abuse of notation, we write  $f(x) = f_x \in \mathbb{R}$  for  $f(x) = f_x \delta_x$ , where  $\delta_x$  denotes a unit mass at  $x$ .

expectation over the value of the contracting problem in the terminal period as a function of posterior beliefs. Let  $P_\mu^e$  denote the expected probability of high output under belief  $\mu$  if the agent exerts effort  $e$ . The optimal contract solves

$$\Pi_t(\mu) = \max_{S,p,w} P_\mu^1 Y + \int_S \left( P_\mu^1 p(s|y_H) + (1 - P_\mu^1) p(s|y_L) \right) \left( \delta \Pi_{t+1}(\mu(s)) - w(s) \right) ds \quad (3)$$

$$\text{s.t. } \int_S \left( P_\mu^1 p(s|y_H) + (1 - P_\mu^1) p(s|y_L) \right) u(w(s)) ds - c \geq U \quad (P)$$

$$\int_S \left( P_\mu^1 p(s|y_H) + (1 - P_\mu^1) p(s|y_L) \right) u(w(s)) ds - c \geq$$

$$\int_S \left( P_\mu^0 p(s|y_H) + (1 - P_\mu^0) p(s|y_L) \right) u(w(s)) ds \quad (\text{IC})$$

$$\int_S p(s|y_H) ds = \int_S p(s|y_L) ds = 1; \quad p(s|y) \geq 0 \quad (S)$$

This is a standard moral hazard problem, but with two added features. First, the principal *chooses an evaluation structure* and the wage cannot be more informative about the agent's output than the evaluation structure it is based on. In particular, the principal can choose to condition the wage on partially informative signals of output instead of output directly. Second, there is a *belief-dependent continuation value*  $\Pi_{t+1}(\mu(s))$ .

Note that the posterior  $\mu(s)$  only realizes if the agent exerted effort. After a deviation, by contrast, the agent makes a different inference. This does not impact the incentive constraint since the principal observes the deviation. The agent does not gain *private* information about his belief and the principal holds him to his reservation value in the second period regardless.<sup>20</sup>

**Proposition 1.** *The optimal contract contains no signals that induce the same belief but are mapped to different wages. The contracting problem can be written as a choice of a distribution over posteriors  $m$  with mean  $\mu$  and support on  $[\underline{\mu}, \bar{\mu}]$ , and a mapping from posteriors to wages.*

While rewriting the choice of a signal structure as a choice of a distribution over posteriors is standard in the literature on Bayesian persuasion, applying this transformation to our contracting problem requires two adaptations. First, note

<sup>20</sup>We solve the problem where this belief-manipulation effect is present in the IC as an extension in Section 5.2.

that the principal designs an information structure about *output*, but the beliefs are about the agent's *type*. Since both spaces are one-dimensional and high quality output is more likely if the agent has a high type, there exists a one-to-one mapping between the two. Second, after a deviation to low effort, the distribution of signals changes. We need to be able to express this change as a function of the posterior distribution. Again, because the mapping from beliefs over output to beliefs over ability is one-to-one, we can find such a transformation (Boleslavsky and Kim, 2017).

Let  $m$  denote the distribution over posteriors and (with slight abuse of notation)  $w$  the mapping from posteriors to utilities associated to  $(S, p, w)$ . It is easy to see that

$$\int_S \left( P_\mu^1 p(s|y_H) + (1 - P_\mu^1) p(s|y_L) \right) \left( \delta \Pi_{t+1}(\mu(s)) - w(s) \right) ds \quad (4)$$

$$= \int m(\hat{\mu}) (\delta \Pi_{t+1}(\hat{\mu}) - w(\hat{\mu})) d\hat{\mu}, \quad (5)$$

and similarly for the participation constraint. To transform the incentive constraint, note first that the original form of the incentive constraint is equivalent to

$$\int_S \left( b + \Delta b \mu \right) \left( p(s|y_H) - p(s|y_L) \right) u(w(s)) ds \geq c \quad (6)$$

An increase in effort increases the probability of high output by  $b + \Delta b \mu$ . This increase affects utility by shifting mass towards signals that are more likely after high output, therefore incentive compatibility requires a sufficiently strong correlation between a signal's responsiveness to high output and the utility delivered after it. Transforming the contracting problem into belief space, it reads

$$\Pi_t(\mu) = \max_{m, w} P_\mu Y + \int m(\hat{\mu}) (\delta \Pi_{t+1}(\hat{\mu}) - w(\hat{\mu})) d\hat{\mu} \quad (7)$$

$$\text{s.t.} \quad \int u(w(\hat{\mu})) m(\hat{\mu}) d\hat{\mu} - c \geq U \quad (\text{P})$$

$$\int (b + \Delta b \mu) \frac{\hat{\mu} - \mu}{(\Delta a + \Delta b) \mu (1 - \mu)} u(w(\hat{\mu})) m(\hat{\mu}) d\hat{\mu} \geq c \quad (\text{IC})$$

$$\int \hat{\mu} m(\hat{\mu}) d\hat{\mu} = \mu; \quad \text{supp}(m) \subset [\underline{\mu}, \bar{\mu}] \quad (\text{BP})$$

The incentive constraint now requires a sufficiently strong correlation between the

*posterior* and utility. This is because signals that are more likely after a good outcome are also associated with a high posterior probability that the agent is the high type. This correlation is rescaled since, depending on the parameters of the problem, this dependence may be more or less strong.

### 3.2 Terminal Period

In the second period, the principal has no continuation value from the relationship. Absent any reason to manipulate the agent’s learning, the only objective in designing the signal structure is to provide incentives cheaply and there is no reason to leave information about output unused. It is optimal to use the most informative signal structure (Grossman and Hart, 1983).

**Proposition 2.** *The optimal contract in the second period uses the fully informative evaluation structure.*

The profit in the terminal period induces a continuation value

$$\int \Pi_2(\hat{\mu})m(\hat{\mu}) d\hat{\mu} \tag{8}$$

for the principal in the first period, where  $m$  is the distribution over posteriors induced by learning from the evaluation in the first-period. We now show that this learning is costly for the principal, since it always reduces her continuation value.

More information about the agent’s ability has two effects. On the one hand, it allows the principal to adapt the contract to the agent’s ability. The contract filters out the nuisance parameter “ability” more effectively and provides incentives for effort more precisely. As a consequence, the wage can be less risky and it is cheaper to provide incentives. This effect is stronger the larger the effect of ability on the probability of high output. On the other hand, the agent also has more information when he decides whether to shirk or exert effort. Consequently, the wage has to be more risky on average in order to satisfy the IC constraint and it is more expensive to provide incentives. In other words, it is easier to satisfy the incentive compatibility constraint in expectation (“ex-ante”) rather than for a more informed agent (“interim”). This effect is stronger the larger the effect of ability on the impact of effort. It dominates and learning reduces profits

whenever complementarities are sufficiently strong, i.e. when the technology is log-supermodular in effort and ability.<sup>21</sup>

**Proposition 3.** *The value of the second period contracting problem,  $\Pi_2$ , is strictly concave in beliefs.*

*Equivalently, consider the expected continuation value induced by distributions over beliefs,  $m, m' \in \Delta([0, 1])$ , where  $m$  is Blackwell less informative than  $m'$ . The principal prefers the less informative distribution*

$$\int \Pi_2(\hat{\mu})m(\hat{\mu}) d\hat{\mu} \geq \int \Pi_2(\hat{\mu})m'(\hat{\mu}) d\hat{\mu}. \quad (9)$$

To see this effect of information on the costs of incentives more concretely, consider the IC constraint in the terminal period,

$$(b + \Delta b\mu)(u(w_H) - u(w_L)) = c. \quad (10)$$

High effort increases the probability of high output by  $P_\mu^1 - P_\mu^0 = b + \mu\Delta b$ . The principal pays a base wage  $w_L$  and adds a bonus  $w_H - w_L$  if and only if output is high (Proposition 2). The utility bonus is inversely proportional to the expected impact of effort,  $b + \mu\Delta b$ , and thus a convex function of the agent's beliefs. Consequently, a greater dispersion of beliefs causes an increase in the expected bonus. The principal wants the agent to stay uninformed, because it is cheaper to pay a bonus that is large enough in expectation than the expected bonus required by an informed agent.

### 3.3 Initial Period

In the first period, our main trade-off is in effect. By Proposition 2, providing incentives for the agent is cheaper in this period if the evaluation structure is more informative, while by Proposition 3 the resulting learning is costly as it increases the expected cost of incentives in the next period.

How is this trade-off resolved in the optimal contract? We employ the tools of information design to characterize the optimal evaluation structure without imposing any exogenous restrictions. While such restrictions, e.g. to a binary

---

<sup>21</sup>This is merely a *sufficient* condition. It is not tight for any nondegenerate utility function. Furthermore, learning is also costly if substitutability is sufficiently strong. A sufficient condition is that the probability of low output is log-supermodular in effort and ability.

evaluation structure, may seem natural in a setting with a binary state and binary output, we know from this literature that they can be with loss of generality. Indeed, since the contracting problem has two constraints – participation and incentive compatibility, results from constrained information design suggest that the optimal evaluation structure may involve up to four signals (Le Treust and Tomala, 2019; Doval and Skreta, 2018).

To analyze this joint information and contract design problem, we make some assumptions on the utility function.

**Assumption 1.** *Let  $w = u^{-1}$  denote the wage function mapping a level of utility to the wage required to provide it. It satisfies*

1. (No incentives at infinity)  $\frac{w(x)}{x} \rightarrow \infty$  as  $x \rightarrow \infty$ .
2. (Bounded changes in curvature)

$$\frac{3(b + \mu\Delta b)\Delta b}{c(a\Delta b - b\Delta a)} \geq \frac{w'''(u_L)}{w''(u_L)} \text{ and } \frac{w'''(u_H)}{w''(u_H)} \geq -\frac{3(b + \mu\Delta b)\Delta b}{c((1-a)\Delta b + b\Delta a)},$$

where  $u_L = U - \frac{a+\mu\Delta a}{b+\mu\Delta b}c$  and  $u_H = U + \frac{1-a-\mu\Delta a}{b+\mu\Delta b}c$ .

3. (Decreasing curvature)  $w''' \leq 0$ .

All three restrictions are sufficient conditions that will be used in the proof of the main theorem. The first condition ensures that an interior solution exists. The principal doesn't find it profitable to provide an arbitrarily high payment with vanishing probability in order to incentivize the agent. The second condition ensures that the shape of the continuation value  $\Pi_2$  is determined unambiguously by the technology and not by changes in the curvature of the utility function. It rules out that the curvature of the utility function changes too quickly. The third condition ensures that the information design problem is governed by the shape of the continuation value. All three conditions are satisfied for CRRA utility ( $u(x) = \frac{x^{1-\gamma}}{1-\gamma}$ ) for  $\gamma \leq \frac{1}{2}$  if the outside utility is sufficiently high.<sup>22</sup> They are always satisfied for  $u(x) = \sqrt{2x}$ .

**Theorem 1.** *Suppose  $u$  satisfies Assumption 1. Then, the optimal evaluation structure in the first period is (essentially) unique. It is binary and tough with*

<sup>22</sup>To see this, note that  $\frac{w'''(x)}{w''(x)} = \frac{2\gamma-1}{1-\gamma} \frac{1}{x}$  for CRRA utility.



$S = \{G, B\}$  and

$$p(G|y_H) = 1 - \sigma, \quad p(B|y_H) = \sigma, \quad p(G|y_L) = 0 \quad p(B|y_L) = 1, \quad (11)$$

for  $\sigma \in [0, 1)$ .

First, the motive to control learning does not increase the complexity of the evaluation structure. While the most informative evaluation is binary, a noisy evaluation can take many forms. The Theorem establishes that the optimal evaluation remains binary. The joint design of wages and information is crucial for this result, it does not necessarily hold when the wage function is fixed exogenously.<sup>23</sup>

Second, the principal uses a noisy binary signal of output as the basis of evaluation. The noise is asymmetric, making the evaluation “tough”: A good evaluation results only if output was high. Low output always results in a bad evaluation, and the bad signal realizes also after high quality output with probability  $\sigma$ . In order to reduce the informativeness of the signal, the principal does not engage in “grade inflation”, but instead measures performance against an “unreasonably” high standard.

The reason for this result is the shape of the continuation value of the principal. While the principal is always information averse, *the degree of information aversion is decreasing in the agent’s posterior* ( $\Pi_2''' > 0$ ). The main objective of the firm is to avoid workers from getting very pessimistic about their ability. To see why, consider again the second period IC,

$$(b + \Delta b\mu)(u(w_H) - u(w_L)) = c. \quad (12)$$

As we discussed previously, the impact of effort,  $b + \Delta b\mu$ , and the required bonus are inversely proportional, which implies that the continuation value is concave. Furthermore, this effect of learning is stronger when the posterior is low. In this case, the agent is pessimistic about the impact of his effort and even a small change in his belief has a large relative effect and causes large changes to the bonus. This leverage effect determines the shape of the continuation value if the curvature of the utility function doesn’t change too much, which is guaranteed by Assumption

---

<sup>23</sup>Georgiadis and Szentes (forthcoming) show that the optimal information structure about effort when there are exogenous costs of acquiring information is binary when wages and the monitoring structure are designed jointly.

1.2. Therefore, the principal's information aversion is larger at low posteriors. In order to raise the low posterior, the optimal monitoring structure pools at the bottom. Since the low evaluation might have been the result of bad luck, it is less damning.

We can provide an interpretable condition for a strictly noisy optimal evaluation in a special case.

**Proposition 4.** *Let  $u(x) = \sqrt{2x}$ . Then evaluation structure is not fully informative ( $\sigma > 0$ ) if and only if*

$$\frac{1}{2} \left( c \frac{\Delta a + \Delta b \mu (1 - \mu) (\bar{\mu} - \underline{\mu})}{b + \Delta b \mu (\bar{\mu} - \mu) (\mu - \underline{\mu})} \right)^2 < \delta \left( \Pi_2(\underline{\mu}) + \Pi_2'(\bar{\mu})(\bar{\mu} - \underline{\mu}) - \Pi_2(\bar{\mu}) \right) \quad (13)$$

The LHS corresponds to the cost of noisier evaluation in period 1, while the RHS is the cost of learning through the continuation value of the principal. A noisy evaluation structure is optimal if and only if the latter cost dominates.

### Proof of Theorem 1

The proof of Theorem 1 poses the challenge of jointly designing an information structure and a wage scheme. Given a wage scheme, the information design problem can be solved by concavification (Aumann and Maschler, 1995; Kamenica and Gentzkow, 2011) taking into account the P and IC constraints (Boleslavsky and Kim, 2017; Le Treust and Tomala, 2019). The constraints make the problem multidimensional so that, although conceptually tractable, concavification is analytically difficult. Conversely, given an information structure, the problem of finding wages is a standard moral hazard problem. This tractable problem provides the starting point for a duality-based approach to such a joint information and incentive design problem, as outlined in Georgiadis and Szentes (forthcoming).

Consider the Lagrangian  $\mathcal{L}$  associated to the contracting problem (7), where we retain (BP) as a constraint, and  $\lambda_P, \lambda_{IC}$  denote the Lagrange multipliers associated

to the participation and incentive constraint, respectively,

$$\begin{aligned} \mathcal{L}(m, w; (\lambda_P, \lambda_{IC})) = & \int \left\{ P_\mu^1 Y + \delta \Pi_2(\hat{\mu}) - w(\hat{\mu}) \right. \\ & + \lambda_P (u(w(\hat{\mu})) - c - U) \\ & \left. + \lambda_{IC} \left( \frac{b + \Delta b \mu}{(\Delta a + \Delta b) \mu (1 - \mu)} (\hat{\mu} - \mu) u(w(\hat{\mu})) - c \right) \right\} d\hat{\mu}. \end{aligned} \quad (14)$$

We will write  $\lambda = (\lambda_P, \lambda_{IC})$  when convenient.

As the proof relies on duality arguments, let us provide a quick summary. The contracting problem is equivalent to

$$\sup_{w, m \text{ s.t. (BP)}} \inf_{\lambda \geq 0} \mathcal{L}(m, w; (\lambda_P, \lambda_{IC})). \quad (15)$$

To see this, note that

$$\inf_{\lambda \geq 0} \mathcal{L}(m, w; \lambda) = \begin{cases} \int m(\hat{\mu}) \left( P_\mu^1 Y + \delta \Pi_2(\hat{\mu}) - w(\hat{\mu}) \right) d\hat{\mu} & \text{if (P)\&(IC) are satisfied} \\ -\infty & \text{else} \end{cases}, \quad (16)$$

the infimum simply wraps the constraints into the objective function. It is always the case that  $\inf \sup \mathcal{L} \geq \sup \inf \mathcal{L}$ , where the supremum is taken over the choice variables and the infimum over the multipliers. If this condition holds with equality, i.e. if we can exchange sup and inf, we say that the optimization problem satisfies *strong duality*.

**Wage Setting** Fix a distribution  $m$  satisfying (BP) and consider the problem of finding optimal wages subject to the participation and incentive constraint. This is a standard moral hazard problem and the optimal the optimal wage schedule follows from pointwise optimization of the Lagrangian. Furthermore, the problem is well behaved, so we have the following Lemma.

**Lemma 1.** *The wage setting problem satisfies strong duality, i.e.*

$$\sup_w \inf_{\lambda \geq 0} \mathcal{L}(m, w; \lambda) = \inf_{\lambda \geq 0} \sup_w \mathcal{L}(m, w; \lambda).$$

**Information Design: Duality** Given Lagrange multipliers  $\lambda$ , the Lagrangian at the optimal wage schedule can be written as an expectation of a function of the

posterior

$$\sup_w \mathcal{L}(m, w; \lambda) = \int \ell^*(\hat{\mu}; \lambda) m(\hat{\mu}) d\hat{\mu} \quad (17)$$

Therefore, the information design problem is of standard form. The principal maximizes the expectation of a function of posteriors,

$$\sup_{m \text{ s.t. (BP)}} \int \ell^*(\hat{\mu}; \lambda) m(\hat{\mu}) d\hat{\mu}, \quad (18)$$

and can therefore be solved via concavification of  $\ell^*$ . But note that this problem takes  $\lambda$  *as given*. This requires another exchange of sup and inf, which can be justified as  $\ell^*$  is continuous and the space of beliefs is compact.

**Lemma 2.** *The information design problem satisfies strong duality, i.e.*

$$\sup_{m \text{ s.t. (BP)}} \inf_{\lambda \geq 0} \int \ell^*(\hat{\mu}; \lambda) m(\hat{\mu}) d\hat{\mu} = \inf_{\lambda \geq 0} \sup_{m \text{ s.t. (BP)}} \int \ell^*(\hat{\mu}; \lambda) m(\hat{\mu}) d\hat{\mu}. \quad (19)$$

**Information Design: Concavification** It remains to solve 18. In order to determine the concavification of  $\ell^*$ , we need to determine its shape as a function of  $\hat{\mu}$ . Using an envelope argument, it is straightforward to show<sup>24</sup> that

$$\begin{aligned} \frac{\partial^2}{\partial \hat{\mu}^2} \ell^*(\hat{\mu}; \lambda) &= \lambda_{IC}^2 \left[ \frac{b + \Delta b \mu}{(\Delta a + \Delta b) \mu (1 - \mu)} \right]^2 \rho'(\lambda_P + \lambda_{IC} \frac{b + \Delta b \mu}{(\Delta a + \Delta b) \mu (1 - \mu)} (\hat{\mu} - \mu)) \\ &\quad + \delta \Pi_2''(\hat{\mu}) \end{aligned} \quad (20)$$

where  $\rho(x) := u(u'^{-1}(\frac{1}{x}))$  denotes the function that translates multipliers and scores to utilities, a function commonly encountered in moral hazard problems. The first term corresponds to the cost of providing incentives in the first period. It is positive, indicating convexity: the principal prefers the most informative evaluation structure in order to reduce agency costs. The second term corresponds to the impact of beliefs on the continuation value. It is negative: the principal wants to keep the agent uninformed in order to reduce agency costs in the next period.

---

<sup>24</sup>In the main text, we suppress boundary conditions related to the non-negativity constraint on wages.

Furthermore, we have that

$$\begin{aligned} \frac{\partial^3}{\partial \hat{\mu}^3} \ell^*(\hat{\mu}; \lambda) = & \lambda_{IC}^3 \left[ \frac{b + \Delta b \mu}{(\Delta a + \Delta b) \mu (1 - \mu)} \right]^3 \rho''(\lambda_P + \lambda_{IC} \frac{b + \Delta b \mu}{(\Delta a + \Delta b) \mu (1 - \mu)} (\hat{\mu} - \mu)) \\ & + \delta \Pi_2'''(\hat{\mu}) > 0 \end{aligned} \quad (21)$$

The shape of  $\ell^*$  has two components. The first term is the impact of the shape of the utility function. For given Lagrange multipliers, it is cheaper to provide incentives at higher posteriors as the curvature of  $w$  is decreasing (Assumption 1.3) and, equivalently,  $\rho'' > 0$ .<sup>25</sup> The second term is determined by the shape of the continuation value. The principal is less information averse for high posteriors.

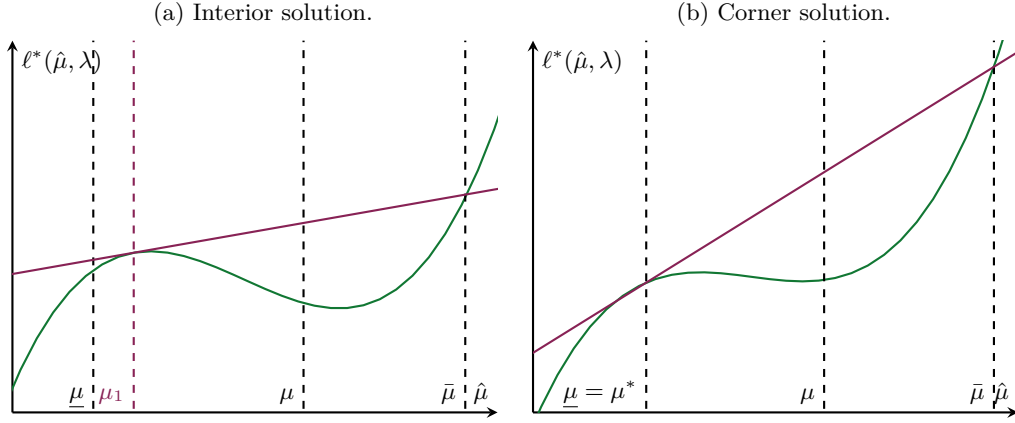


Figure 1: The concavification of  $\ell^*$  at  $\mu$ .

There are three possible cases. If  $\lambda_{IC}$  is sufficiently small, the objective  $\ell^*$  is strictly concave and the optimal information structure is uninformative. Clearly, this cannot be the case in the solution of (7), since the incentive constraint cannot be satisfied without any information. As  $\lambda_{IC}$  increases, we reach a region where the optimal information structure is partially revealing. As  $\ell^*$  is convex for high posteriors in that case, it is fully informative at the top and uses partial pooling at the bottom (Fig. 1a). Finally, as  $\lambda_{IC}$  increases further, the costs of incentives

<sup>25</sup>Note that this effect is purely “partial optimality”: in the solution, the multipliers have to adjust to make sure that wage dispersion is sufficient to satisfy incentive compatibility. I conjecture that the restriction to utility functions with  $\rho'' > 0$  is far from necessary: Instead, it is a result of the proof approach that requires solving the problem for all multipliers and relies on properties of the solution that are uniform across multipliers.  $\frac{\partial^3}{\partial \hat{\mu}^3} \ell^*(\hat{\mu}; \lambda)$  is positive for all multipliers only if  $\rho'' \geq 0$ .

overwhelm the gains from concealing information and the evaluation structure is fully informative (Fig. 1b).

**Lemma 3.** *For any  $\lambda$ , the optimal evaluation structure is unique and induces at most two posteriors. It induces the highest feasible posterior  $\bar{\mu}$  with probability  $m(\bar{\mu}) \in [0, \frac{\mu - \underline{\mu}}{\bar{\mu} - \underline{\mu}}]$  and a low posterior,  $\mu^* \in [\underline{\mu}, \mu]$  with  $m(\mu^*) \in [\frac{\bar{\mu} - \mu}{\bar{\mu} - \underline{\mu}}, 1]$ .*

**The Simplified Problem** We can simplify the general problem (7) using the properties of optimal evaluation structures from the previous lemma, i.e. we can restrict attention to binary information structures where the good signal only realizes after high output. This simplified problem is

$$\max_{\mu^*, m^*, w_l, w_h} P_\mu^1 Y + m^* [\delta \Pi_2(\mu^*) - w_l] + (1 - m^*) [\delta \Pi_2(\bar{\mu}) - w_h] \quad (22)$$

$$\text{s.t. } m^* u(w_l) + (1 - m^*) u(w_h) - c \geq U \quad (\text{P}^S)$$

$$\frac{b + \Delta b \mu}{(\Delta a + \Delta b) \mu (1 - \mu)} [m^* (\mu^* - \mu) u(w_l) + (1 - m^*) (\bar{\mu} - \mu) u(w_h)] \geq c \quad (\text{IC}^S)$$

$$m^* \mu^* + (1 - m^*) \bar{\mu} = \mu; \quad \mu^* \in [\underline{\mu}, \mu] \quad (\text{BP}^S)$$

where  $\mu^*$  denotes the posterior after the bad evaluation,  $m^*$  denotes the probability of a bad evaluation and  $w_l, w_h$  denote the low and high wage, respectively.

**Lemma 4.** *The simplified contracting problem (22) has a unique solution. The optimal information structure is non-degenerate ( $\mu^* < \mu$ ).*

The condition for an interior solution is derived in the Appendix. QED.

### 3.4 Analysis of the Solution and Comparative Statics

The properties of an interior solution are pinned down by tangency condition

$$\ell^*(\mu^*, \lambda(\mu^*)) + \frac{\partial \ell^*}{\partial \hat{\mu}} \Big|_{(\hat{\mu}, \lambda) = (\mu^*, \lambda(\mu^*))} (\bar{\mu} - \mu^*) = \ell(\bar{\mu}, \lambda(\mu^*)), \quad (23)$$

determining the posterior  $\mu^*$  in the concavification of  $\ell^*$  (Fig. 1). Note, however, that this condition does not correspond to the concavification of a given function. Instead, there is an additional dependence on  $\lambda(\mu^*)$ . This term is present because we are not solving an information design problem given payoffs, but design payoffs

and information jointly, subject to a participation and and incentive compatibility constraint. For the graphical representation of our analysis this implies that, as we vary the tangent point in the figure to find the optimal  $\mu^*$ , not only the tangent line but the whole function  $\ell^*$  shifts.

Under the assumption that  $u(x) = \sqrt{2x}$  we can transform (23) into a more interpretable form:

$$\frac{c^2}{2} \left( \frac{(\Delta a + \Delta b)\mu(1 - \mu)}{b + \Delta b\mu} \frac{\bar{\mu} - \mu^*}{(\bar{\mu} - \mu)(\mu - \mu^*)} \right)^2 = \delta (\Pi_2(\mu^*) + \Pi_2'(\bar{\mu})(\bar{\mu} - \mu^*) - \Pi_2'(\bar{\mu})) \quad (24)$$

The LHS is the benefit from a more informative evaluation structure in period one. A more precise signal about output decreases agency costs today. This effect is larger if agency costs ( $\frac{c}{b + \Delta b\mu}$ ) are already high and if a large dispersion of posteriors is required for a given level of information about output (since output is very informative,  $\Delta b\mu(1 - \mu)$  large). The RHS is the cost of a more informative information structure through learning. A more precise signal today allows learning and thereby increases average agency costs in the next period. Indeed, the RHS is a measure of the concavity of the continuation value.

The optimal degree of shrouding,  $\sigma$ , is pinned down by the lower posterior belief  $\mu^*$  according to

$$\sigma(\mu^*) = \frac{1 - P_\mu \mu^* - \underline{\mu}}{P_\mu \bar{\mu} - \mu^*} \in [0, 1], \quad (25)$$

which follows from inverting Bayes rule. It is increasing in  $\mu^*$ ; if the principal wants to cushion bad news, she needs to pool more on the bad signal.

**Proposition 5.** *Suppose  $u(x) = \sqrt{2x}$ . The optimal level of shrouding  $\sigma$  is*

- (1) *weakly increasing in the discount factor  $\delta$*
- (2) *weakly decreasing in the costs of effort in the first period, weakly increasing in the costs of effort in the second period, and independent of a common increase in the cost of effort.*

*All comparisons are strict at interior  $\sigma$ .*

Both comparative statics illustrate the trade-off between the cost of incentives in the first and second period. As the second period becomes more important, the

evaluation structure becomes less informative. Higher costs of effort in the first period make economizing on agency costs in that period more important, thus the evaluation structure becomes more informative.

## 4 Preserving and Correcting Misperceptions

So far, we have argued that noisy and tough evaluation is the optimal way to preserve uncertainty about the agent’s ability while providing incentives. We assumed that the principal and the agent agree about the situation, i.e. that they share a common prior. There is some evidence suggesting, however, that beliefs concerning the impact of effort on outcomes – which are the driving factor of our results – may be systematically biased. Overestimation of one’s abilities has been demonstrated in several laboratory contexts as well as in the workplace.<sup>26</sup> Overconfident workers overestimate their type and hence, given the complementarity between effort and ability, the importance of their contribution. Other biases can also affect beliefs about the impact of effort, for example the illusion of control, a tendency to overestimate the impact of individual choices on outcomes that also depend on chance, or the belief in a “just world”.<sup>27</sup> Some individuals are also systematically underconfident and this trait is common in some groups.<sup>28</sup>

---

<sup>26</sup>See, for example [Larwood and Whittaker \(1977\)](#) for early evidence that individuals overestimate their abilities in a laboratory setting, ([Burks et al., 2013](#)) for a more recent incentivized study. Overconfidence is also present in tournaments ([Park and Santos-Pinto, 2010](#)) and among store managers ([Huffman et al., 2019](#)).

<sup>27</sup>[Langer \(1975\)](#) defines the illusion of control broadly as "an expectancy of a personal success probability inappropriately higher than the objective probability would warrant". The typical experiment establishes increased optimism about the outcome of a lottery in situations involving "choice, stimulus or response familiarity, passive or active involvement or competition". The fact that most experiments involve pure chance is intended as an extreme condition, suggesting that "the effects should be far greater when they are introduced into situations when there already is an element of control". But note [Charness and Gneezy \(2010\)](#); [Filippin and Crosetto \(2016\)](#), who find no evidence of illusion of control in two main experimental paradigms with monetary incentives.

According to just-world belief, effort and more generally good deeds are rewarded in the world. Such attitudes vary widely across countries and appear at best weakly related to true level of meritocracy. See [Lerner \(1980\)](#), and [Bénabou and Tirole \(2006\)](#) and the references therein for a discussion of the evidence.

<sup>28</sup>There is some evidence that women tend to be underconfident, for example ([Niederle and Vesterlund, 2007](#); [Hügelschäfer and Achtziger, 2014](#)).



In this section, we analyze the optimal contract when the agent is not merely uncertain about his type, but enters the relationship with a systematic misperception. The principal now has an additional motive to shape learning, namely to affect the average posterior of the agent. An agent who overestimates his ability is more profitable because he is easier to incentivize. The principal would like to preserve this profitable misconception. Is this still achieved via tough evaluations or does she use a lenient information structure, akin to grade inflation, as the optimal way to preserve optimism?<sup>29</sup>

#### 4.1 Contracting with Heterogeneous Priors

We solve the contracting problem with heterogeneous priors. The agent again has a prior belief  $\mu$  that he has high ability. The principal, by contrast, has a prior belief  $\eta \in \{0, 1\}$ .<sup>30</sup> When  $\eta = 0$ , the principal is sure that the agent has low ability and we say that the agent is overconfident. When  $\eta = 1$ , by contrast, the principal is sure the agent has high ability and the agent is underconfident. The two players agree to disagree and update their priors using Bayes rule.<sup>31</sup>

We maintain our restrictions on the technology, namely that effort is productive ( $b \geq 0$ ), the high type is more productive ( $\Delta a \geq 0$ ), the technology is log-supermodular ( $a\Delta b - b\Delta a > 0$ ) and that the outside option is sufficiently attractive to ensure an interior solution ( $U > \frac{a+b}{b}c$ ). To focus on the effect of heterogeneous beliefs on the problem, we assume that  $u(x) = \sqrt{2x}$ .<sup>32</sup>

#### The Transformation to Belief Space

As before, we will transform the contracting problem and write it as the choice of a distribution of posterior beliefs and a wage function. However, the principal

---

<sup>29</sup>Indeed, supporting students' self-esteem is often cited as a reason for grade inflation in schools and universities (Boretz, 2004).

<sup>30</sup>Generally, it is reasonable to believe that the principal has better knowledge than the agent about his match-specific ability. The assumption that the principal is certain about the agent's match-specific ability is crucial for tractability in the case of heterogeneous priors. This is because the continuation value now depends both on the agents belief and the level of disagreement. With either identical priors or one degenerate prior, these two variables are simple. If these restrictions don't hold, the problem can still be rewritten in one dimension, but the information design problem is not tractable.

<sup>31</sup>While it would be interesting to analyze the design problem with non-Bayesian players, there is little work on information design tools for such a setting.

<sup>32</sup>This choice of utility function ensures that the curvature of the cost of wages,  $w = u^{-1}$ , (as a function of utility) is constant in the principal's problem.

and the agent now have heterogeneous priors. In particular, the principal does not learn and therefore the relevant posteriors are those of the agent. In addition, the principal and the agent have different beliefs over the induced distribution of these posteriors. Let  $m$  denote the distribution according to the agent's belief and  $m_P$  this distribution according to the principal. The distribution over posteriors satisfies Bayes plausibility according to the agent,

$$\int \hat{\mu} m(\hat{\mu}) d\hat{\mu} = \mu \quad (26)$$

but, generically, not according to the principal. We can write the distribution over posteriors under the principal's prior belief,  $m_P$ , as a transformation of  $m$ , as follows. Let  $s$  be the signal inducing posterior  $\hat{\mu}(s)$ .<sup>33</sup> Then, the probability of  $\hat{\mu}(s)$  according to the agent is

$$m(\hat{\mu}(s)) = p(s|y_L) + (a + b + (\Delta a + \Delta b)\mu) [p(s|y_H) - p(s|y_L)] \quad (27)$$

According to the principal, this event has probability

$$\begin{aligned} m_P(\hat{\mu}(s)) &= p(s|y_L) + (a + b + (\Delta a + \Delta b)\eta) [p(s|y_H) - p(s|y_L)] \\ &= m(\hat{\mu}(s)) + (\eta - \mu)(\Delta a + \Delta b) [p(s|y_H) - p(s|y_L)] \\ &= \left[ \eta \frac{\hat{\mu}(s)}{\mu} + (1 - \eta) \frac{1 - \hat{\mu}(s)}{1 - \mu} \right] m(\hat{\mu}(s)) \end{aligned} \quad (28)$$

Hence, we can follow the approach of [Alonso and Câmara \(2016\)](#) to Bayesian persuasion with heterogeneous priors and solve for the distribution  $m$  while the transformation factor  $D^\eta(\mu, \hat{\mu}) := \left[ \eta \frac{\hat{\mu}}{\mu} + (1 - \eta) \frac{1 - \hat{\mu}}{1 - \mu} \right]$  takes the heterogeneous priors into account.<sup>34</sup>

---

<sup>33</sup>This signal is unique without loss of generality by a straightforward extension of Proposition 1.

<sup>34</sup>Note that, in contrast to [Alonso and Câmara \(2016\)](#), the priors of the principal and the agent on the state space are *not* mutually absolutely continuous. The transformation method (as opposed to information design with surprises, [Galperti, 2019](#)) is still applicable since the posterior needs to be measurable with respect to a noisy signal of the state, namely output. This restriction keeps the belief transformation bounded. To be more precise, in our framework the principal designs an information structure about output, which implies a posterior about the type. Beliefs about the distribution of output are heterogeneous, but mutually absolutely continuous. There is a 1:1 mapping from beliefs about output to posteriors over the type.

The contracting problem with heterogeneous beliefs is thus

$$\Pi_t^\eta(\mu) = \max_{m,w} P_\eta^1 Y + \int (\delta \Pi_{t+1}^\eta(\hat{\mu}) - w(\hat{\mu})) D^\eta(\mu, \hat{\mu}) m(\hat{\mu}) d\hat{\mu} \quad (29)$$

$$\text{s.t. } \int u(w(\hat{\mu})) m(\hat{\mu}) d\hat{\mu} - c \geq U \quad (\text{P})$$

$$\int (b + \Delta b \mu) \frac{\hat{\mu} - \mu}{\Delta b \mu (1 - \mu)} u(w(\hat{\mu})) m(\hat{\mu}) d\hat{\mu} \geq c \quad (\text{IC})$$

$$\int \hat{\mu} m(\hat{\mu}) d\hat{\mu} = \mu; \quad \text{supp}(m) \subset [\underline{\mu}, \bar{\mu}] \quad (\text{BP})$$

## 4.2 Terminal Period

Our results about the problem in the final period extend to the setting with heterogeneous priors. There is no reason to shape learning, so the principal prefers as much information as possible to incentivize effort as cheaply as possible. Therefore, the optimal evaluation structure in the final period is fully informative.<sup>35</sup>

In order to evaluate the impact of learning on the continuation value of the principal, we need to take into account the measure transform and consider

$$D^\eta(\mu, \hat{\mu}) \Pi_2^\eta(\hat{\mu}) \quad (30)$$

Learning creates a dispersion of the agent's posterior. This is costly for the principal, since  $\frac{\partial^2}{\partial \hat{\mu}^2} \Pi_2^\eta(\hat{\mu}) < 0$  for the reasons discussed in the previous section. In addition, learning now also affects the expected posterior under the principal's belief. This *drift* has two effects. First, the disagreement between the principal and the agent decreases. This makes it harder to gamble on their belief difference and reduces profits. Second, the agent move towards the truth on average. Since gambling is limited due to risk aversion, this second effect dominates. If the agent is optimistic about the impact of effort ( $\eta = 0$ ), this means he becomes less optimistic as he learns. Since optimism is profitable, therefore the principal has an additional incentive to sabotage learning. If, instead, the agent is pessimistic ( $\eta = 1$ ), he becomes less pessimistic on average, which is good for the principal.

The total effect of learning combines the two forces of increased dispersion and drift. Therefore, learning reduces profits with optimism and has an ambiguous impact with pessimism. Let us summarize the preceding discussion.

---

<sup>35</sup>In contrast to Proposition 2, this result does not follow readily from standard results about information in moral hazard problems, but requires an extension of the usual argument.

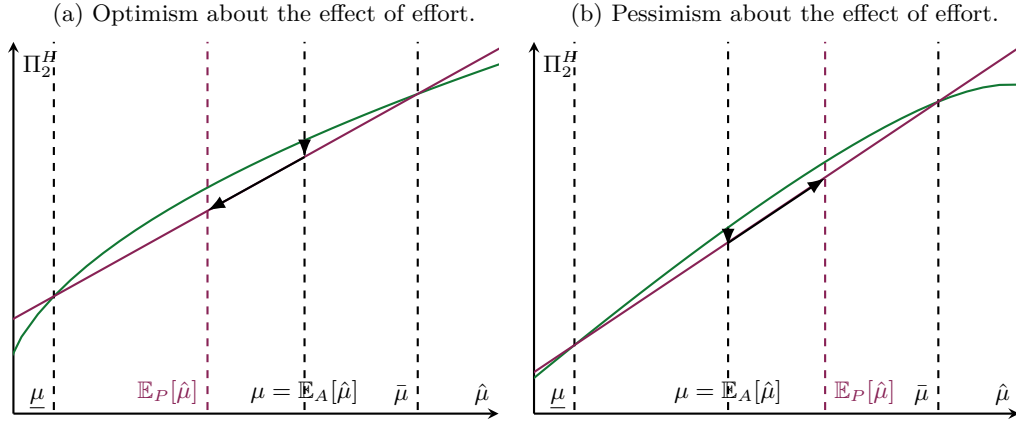


Figure 2: The effect of information on the principal's continuation value.

**Proposition 6.** *Consider the contracting problem with heterogeneous beliefs. In the terminal period, the optimal evaluation structure is fully informative. The value of the second period contracting problem,  $\Pi_2^\eta$ , is strictly increasing and concave in the agent's posterior.*

The impact of information is determined by

$$\frac{\partial^2}{\partial \hat{\mu}^2} (D^\eta(\mu, \hat{\mu}) \Pi_2^\eta(\hat{\mu})) = \underbrace{\left[ \eta \frac{\hat{\mu}}{\mu} + (1 - \eta) \frac{1 - \hat{\mu}}{1 - \mu} \right] \Pi_2^{\eta''}(\hat{\mu})}_{\text{dispersion}} + \underbrace{2 \frac{\eta - \mu}{(1 - \mu)\mu} \Pi_2^{\eta'}(\hat{\mu})}_{\text{drift}} \quad (31)$$

It is negative if the agent is overconfident ( $\eta = 0$ ). If the agent is underconfident ( $\eta = 1$ ) the sign is ambiguous. In particular,

- if  $b > 0$ , there exists a threshold  $\bar{U}$  such that the principal is information loving if  $U \geq \bar{U}$ , and
- if  $a(1 - a) > (4a - 1)(\Delta a + \Delta b)$ , there exists a threshold  $\bar{b} > 0$  such that the principal is information averse if  $b < \bar{b}$ .

For high outside utility, the drift effect dominates and the principal prefers to reveal information in order to eliminate costly underconfidence. If the baseline effectiveness of effort is low ( $b$  is small), the impact of beliefs on the required bonus is highly levered, the dispersion effect dominates and the principal prefers to slow learning.

## Initial Period

The shape of the optimal evaluation structure is determined by two factors: First, based on to continuation value, is the principal information averse and how does this information aversion change as a function of the posterior? This effect is similar to the common prior case, with the addition of the impact of the drift effect. Second, how do the costs of delivering utility change as a function of the posterior? This effect is not present with common priors and stems directly from the heterogeneity of beliefs.

Let us start with the familiar first effect. Taking into account the measure transform, the change in information aversion is determined by

$$\frac{\partial^3}{\partial \hat{\mu}^3} (D^\eta(\mu, \hat{\mu})\Pi_2^\eta(\hat{\mu})) = \underbrace{\left[ \eta \frac{\hat{\mu}}{\mu} + (1 - \eta) \frac{1 - \hat{\mu}}{1 - \mu} \right] \Pi_2^{\eta''' }(\hat{\mu})}_{\text{dispersion}} + \underbrace{3 \frac{\eta - \mu}{(1 - \mu)\mu} \Pi_2^{\eta''' }(\hat{\mu})}_{\text{drift}} \quad (32)$$

With an overconfident agent, both effects go in the same direction: A dispersion in beliefs has a higher leverage and is therefore more costly if the agent thinks effort is not very effective, i.e. at low posteriors. Similarly, the impact if the drift is stronger and therefore more costly at low posteriors, since a given decrease in the expected impact of effort has a higher leverage. Therefore, the information aversion of the principal decreases as she induces higher posteriors.

With an underconfident agent, there is a trade-off. A dispersion in beliefs again has a higher leverage and is therefore more costly if the agent thinks effort is not very effective, i.e. at low posteriors. Therefore, inducing dispersion at low posteriors is more costly. Also the drift effect is also stronger with higher leverage at low posteriors, but since the drift effect is desirable with a pessimistic agent, this means that inducing drift at a low posterior is more profitable. Therefore, the total effect is ambiguous, the dispersion effect pushes towards increasing information aversion, while the drift effect pushes towards decreasing information aversion.

Let us turn to the direct effect of belief heterogeneity. It is cheaper for the principal to provide utility to the agent in states that the agent believes to be more likely than the principal. For an overconfident agent, that is a high posteriors, for an underconfident agent, that is at low posteriors.

**Theorem 2.** *The optimal evaluation structure in the first period is unique (up to renaming), binary and uses partial pooling. Let  $S = \{G, B\}$  denote the signal space and  $\sigma \in [0, 1)$  the shrouding parameter*

- *If the agent is overconfident ( $\eta = 0$ ), the optimal evaluation structure is (weakly) strict, i.e.*

$$p(G|y_H) = 1 - \sigma, \quad p(B|y_H) = \sigma, \quad p(G|y_L) = 0 \quad p(B|y_L) = 1. \quad (33)$$

- *If the agent is underconfident ( $\eta = 1$ ), the optimal evaluation structure is (weakly) lenient, i.e.*

$$p(G|y_H) = 1, \quad p(B|y_H) = 0, \quad p(G|y_L) = \sigma \quad p(B|y_L) = 1 - \sigma. \quad (34)$$

Inducing high posteriors is always appealing for overconfident agents. All three effects work in the same direction. For underconfident agents, the drift effect and the direct effect of heterogeneous priors together are strong enough to jointly overpower the increased cost of dispersion associated to providing information at low posteriors. To realize separation at the top and pooling at the bottom (in posterior space), the evaluation is lenient.

## 5 Discussion and Extensions

In the previous sections, we made several assumptions to ensure that the agent's posterior belief is the only state variable of the problem and that no party can acquire endogenous private information. Now, we relax those assumptions and discuss the impact on our results.

### 5.1 Private Information Acquisition

In some settings, it may be possible for the firm to privately observe additional information about the worker's output without disclosing it or using it as a basis of wages in the same period. We analyze this case and show that there exist natural equilibria that replicate the optimal contract. Furthermore, if the firm can commit not to acquire private information, it will.

Consider the model with symmetric priors. For simplicity of exposition, we assume that the principal uses a fully informative evaluation in the second period, as shown to be optimal in Proposition 2.<sup>36</sup> In the first period, the principal now also designs a *private* evaluation structure. Neither this information structure nor its realizations are observed by the agent, and we allow its distribution to depend on the realization of the public signal. Writing the problem in belief space, the principal designs a joint distribution of agent and principal posteriors  $m_P(\mu_P, \hat{\mu})$ , with  $\text{supp}(m_P) \subset [\underline{\mu}, \bar{\mu}]^2$ . The marginal on the agent's posterior,  $m(\hat{\mu}) = \int m_p(\mu_P, \hat{\mu}) d\mu_P$ , is observed by the agent. The distribution satisfies Bayes plausibility for both players,  $\int \hat{\mu} m(\hat{\mu}) d\hat{\mu} = \mu$  and  $\int \mu_P m_p(\mu_P, \hat{\mu}) d\mu_P = \hat{\mu}$ .

The two-period contracting problem now induces a dynamic game with incomplete information. A perfect Bayesian equilibrium consists of (1) an evaluation structure  $m_P$  satisfying the above conditions, (2) a wage function  $w : \hat{\mu} \rightarrow w(\hat{\mu}) \in \mathbb{R}_+$ , (3) a first-period strategy of the worker mapping the evaluation and wage scheme to participation and effort choices,  $m, w \rightarrow \{0, 1\}^2$ , (4) a second period contract offer,  $(\mu_P, \hat{\mu}) \rightarrow (w_L, w_H)(\mu_P, \hat{\mu})$ , (5) a belief system for the agent over his type and the information structure chosen by the principal, as a function of the posterior and the contract offer,  $(\hat{\mu}, w_L, w_H) \rightarrow \Delta([0, 1] \times \Delta[0, 1]^2)$ , satisfying sequential rationality and consistency.

The outcome of Theorem 1 is achieved as the unique equilibrium in a natural class. A PBE is said to have *passive beliefs* if the second period belief of the agent is independent of the contract offer and equal to the posterior  $\hat{\mu}$  induced by the first period signal.<sup>37</sup>

*Remark 1.* The (essentially unique) equilibrium with passive beliefs is outcome-equivalent to the optimal contract characterized in Theorem 1. This equilibrium is principal preferred among all PBE of the game.

The intuition for this result is simple. For the principal facing an agent with passive belief  $\hat{\mu}$ , the optimal contract in the second period satisfies both P and IC with equality. Therefore, the private information of the principal is of no use, and

---

<sup>36</sup>This restriction is without loss on path, as a fully informative evaluation structure remains optimal for the principal. Off path, the restriction reduces the degrees of freedom for deviations, but the equilibrium can be extended naturally.

<sup>37</sup>Orlov et al. (forthcoming) assume passive beliefs to show that the solution to their dynamic persuasion problem is robust to exogenous private information of the sender. Passive beliefs are also a common assumption in games with unobserved bilateral contracts, e.g. Hart and Tirole (1990); Brunnermeier and Oehmke (2013).

the continuation value induced on the first period is the same as in Propositions 2 and 3. Consequently, the principal’s choice of  $m_P$  is equivalent to the first-period problem.<sup>38, 39</sup> To see that this equilibrium is principal preferred, note that in any PBE both the participation constraint and the incentive compatibility constraint need to be satisfied on the equilibrium path. The optimal contract is the best contract satisfying these restrictions. Any information used and thereby revealed in the second period could have been revealed in the first period, thus reducing agency costs.

The principal prefers to commit not to reveal the information in the second period. Passive beliefs provide such a form of commitment. Similarly, consider the game when the principal’s choice of information structure  $m_P$  – both for private and public signals – is observed. Then, we have the following.

*Remark 2.* When the information structure is observed, any equilibrium<sup>40</sup> is outcome equivalent to the optimal contract characterized in Theorem 1.

## 5.2 Unobservable Effort

In the main sections, we assume that effort is observed but not contractible. This ensures that even after a deviation, the principal and the agent have common knowledge about their beliefs over the agent’s type.

Assume instead that effort is not observed by the principal. This does not affect beliefs on equilibrium path, since the conjectured effort is correct. After a

---

<sup>38</sup>Common refinements for signaling games, such as the intuitive criterion (Cho and Kreps, 1987) or D1 (Cho and Sobel, 1990), do not apply as they require the set of types of the principal to be fixed, which is not the case in our game. There are also no proper subgames to which they could be applied. Ekmecki and Kos (2019) analyze a signaling game when the sender chooses whether to acquire full information about his binary type or not, applying a form of never weak best response. Generalizing this kind of analysis to this extension is left for future research.

<sup>39</sup>If we nevertheless apply the reasoning of the intuitive criterion loosely to the contract offer game in the second period, it does not satisfy the requirement. This is because the principal’s types with posteriors above those of the agent have a deviation that allows them to separate. This deviation, however, may not be the most intuitive psychologically. Compared to the pooling contract, the new contract features a lower bonus and delivers lower utility to the agent both under the original and under any plausible posterior belief. One may conjecture that workers see such a contract offer less as a gesture of trust – as the intuitive criterion requires – but as a slight that demonstrate that the principal does not value their continued employment.

<sup>40</sup>Among PBE which satisfy the following natural restriction, a form of no-signaling-what-you-don’t-know: After observing the information structure  $m_P$  and signal  $\hat{\mu}$ , his belief is always supported on the convex hull of the support of  $m_P(\cdot, \hat{\mu})$ .



deviation to  $e_t = 0$ , however, the agent updates his beliefs according to

$$\tilde{\mu}(s) = \mu \frac{p(s|y_L) + (a + \Delta a) [p(s|y_H) - p(s|y_L)]}{p(s|y_L) + (a + \mu \Delta a) [p(s|y_H) - p(s|y_L)]} \quad (35)$$

while the principal continues to use the on-path updating rule (1). Hence, depending on the signal realization, the agent will be less (resp. more) optimistic about his type in the second period and the contract offered by the principal will violate (resp. over-satisfy) the incentive compatibility constraint.<sup>41</sup> A deviation in the first period is more profitable for the agent because of this belief-manipulation effect.<sup>42</sup>

Let us now analyze this model, assuming that  $\Delta a = 0$ . This condition ensures that the agent does not learn about his type after a deviation and simplifies the problem considerably. Note that the problem in the second period is unchanged: The modification only affects continuation beliefs, which are irrelevant in the terminal period. In the first period, we need to modify the incentive-compatibility constraint in order to take the belief-manipulation effect into account.

Let  $w_L(\hat{\mu}(s))$  denote the optimal low wage in the second period problem with belief  $\hat{\mu}(s)$ . The IC reads

$$\int_S (p(s|y_L) + (a + b + \mu \Delta b) [p(s|y_H) - p(s|y_L)]) [w(s) + U] ds - c \geq \int_S (p(s|y_L) + a [p(s|y_H) - p(s|y_L)]) \cdot \left[ w(s) + \max \left\{ w_L(\hat{\mu}(s)) + P_\mu^1 \frac{c}{b + \Delta b \hat{\mu}(s)} - c, w_L(\hat{\mu}(s)) + P_\mu^0 \frac{c}{b + \Delta b \hat{\mu}(s)} \right\} \right] ds \quad (36)$$

This condition is now dynamic: If the agent does not deviate (first line), he will obtain his reservation utility  $U$  in the final period. If effort were observable, this

<sup>41</sup>This assumes that the principal does not elicit the agent's belief at the beginning of the second period. For truth telling to be incentive compatible, it would need to be preferable to imitating the type that realizes on path, however. Hence, a screening mechanism in the second period cannot reduce the post-deviation payoff and therefore does not affect the optimal contract.

<sup>42</sup>This effect is central in the analysis of many models of moral hazard with learning, e.g. [Prat and Jovanovic \(2014\)](#); [Demarzo and Sannikov \(2017\)](#). [Bhaskar and Mailath \(2019\)](#) show that this motive implies that the costs of providing incentives using spot contracts grows unboundedly with the length of the time horizon in a model similar to ours, but with learning from output. It is doubtful whether the design of the information structure can reverse this conclusion and we conjecture that implementing high effort does not remain profitable for a long horizon with unobservable effort in our model.

would also be the case after a deviation, so this term would cancel. Since effort is not observable, he acquires private information about his type after a deviation and has a nontrivial choice in the second period between exerting effort (the first term of the max) and shirking (the second term in the max). The former is optimal if he is more optimistic after the deviation ( $\mu > \hat{\mu}(s)$ ): The principal believes that the signal that realized is indicative of a low type and offers a correspondingly high bonus in the next period. The agent exerts effort and experiences a net gain. The latter is optimal if he is more pessimistic after the deviation ( $\mu < \hat{\mu}(s)$ ): The principal believes that the signal that realized is indicative of a high type and offers a correspondingly low bonus in the next period. The agent does not exert effort and thereby receives his reservation utility, avoiding the loss from the low bonus. Since the agent can reap the gain and avoid the loss, acquiring private information renders a deviation from high effort more profitable.

We can translate this dynamic IC into belief space and write

$$\int \left\{ \frac{(b + \mu\Delta b)}{\mu(1 - \mu)\Delta b} (\hat{\mu} - \mu) u(w(\hat{\mu})) - \left[ 1 - \frac{(b + \mu\Delta b)}{\mu(1 - \mu)\Delta b} (\hat{\mu} - \mu) \right] \max\{0, c\Delta b \frac{\mu - \hat{\mu}}{b + \hat{\mu}\Delta b}\} \right\} m(\hat{\mu}) d\hat{\mu} \geq c \quad (37)$$

Transformed in this fashion, the problem is amenable to an analysis along the lines of Theorem 1. The added complexity, however, is that kink in the incentive compatibility constraint introduces a kink in the Lagrangian of the problem.

*Remark 3.* The Lagrangian of the first period problem is concave-convex, with a concave kink at the prior,  $\hat{\mu} = \mu$ . The optimal evaluation structure therefore consists of

1. a high signal that realizes only if output was good and results in the highest feasible posterior  $\bar{\mu}$ ,
2. (possibly) a neutral signal that results in an unchanged posterior  $\mu$ ,
3. a low signal associated with posterior  $\mu^* \in [\underline{\mu}, \mu)$ .

Conditional on an informative realization from the evaluation, the signal structure is as before. The kink in the IC constraint, however, raises the possibility of a third, uninformative signal. This is because the low signal realization is costly, as it invites belief-manipulation (37). In numerical simulations however, this possibility was never realized and we conjecture that the neutral signal is never part of the optimal contract.

### 5.3 Long-Run Commitment

In the main sections, we assumed that the principal does not have commitment across periods. This is not crucial for our results. What is crucial, however, is that the principal cannot backload all information.

To see this, suppose that the principal can commit to wages that depend on output in both periods and are revealed and paid at the end of the employment relationship. Then, informative wages do not lead to learning and hence using a fully informative evaluation is optimal. Indeed, if the principal can decide when to reveal wages, backloading information in this fashion is optimal. Common impediments against the backloading of payments other than period-by-period contracting can also break this result. Suppose for instance that the agent is less patient than the principal. In the extreme case of a myopic agent, only the current payments of the principal matter for payoffs and the period-by-period contracting problem is equivalent to the one discussed in the main sections. Noisy and (weakly) tough performance evaluation is again optimal.<sup>43</sup>

Our results continue to hold if the principal can postpone payments, but not information. Suppose that the first period contract specifies not only a wage this period, but also a continuation value.<sup>44</sup> Our results generalize to this model; the optimal evaluation structure can be noisy and is (weakly) tough.

Our results can be extended to the case where the principal can engage in partial informational backloading, as follows: The agent does not observe the wages he receives. Instead, there is a probability  $\alpha \in [0, 1]$  that he observes the outcome of his evaluation (e.g. the agent overhears the management talking about it). Importantly, the firm observes whether the agent did observe the evaluation or not. The utility specification and timing of contracting remains the same as in the main sections. The optimal contract in the first period is equivalent to the

---

<sup>43</sup>A continuity argument suggests that the result generalizes to interior impatience. A full analysis of this problem is left for future research.

<sup>44</sup>Note that specifying only a continuation value but not the exact way this value is delivered is potentially with loss in this setting. The outcome with commitment to a continuation value can be replicated by a two-period contract only if the principal can condition on the true posterior of the agent. This, however, would render effort contractible. For a model of moral hazard with long-run contracting, we also need to assume unobservable effort. Then, however, the continuation value is not a sufficient statistic for the influence of future payments on the current contract. In particular, the second period IC may be slack in the optimal contract after some realizations of the first-period signal.

solution to our problem with discount factor  $\alpha\delta$ . A lower chance of discovery  $\alpha$  leads to a more informative performance evaluation, lower agency costs.

## 6 Concluding Remarks

Our model demonstrates why it can be in a principal's interest to base incentives on a noisy evaluation of the agent's performance, even when the principal could measure true output and commit to contingent wages. The underlying insight is that output contains information both about effort, which she wants to ascertain and incentivize, and the agent's match-specific ability, which she would like to keep shrouded.

The optimal performance evaluation is *tough*: Good performance is not always recognized, but bad performance is always punished. Such tough evaluation ensures that even after a bad evaluation, the agent is not too pessimistic about his type. This is optimal because learning is especially costly at low posteriors, as a given change in beliefs has a large impact relative to the small expected efficiency of effort. One way a firm can commit to tough evaluation is through the selection and training of evaluators. Unreasonably strict supervisors and drill-sergeant mentality is part of the optimal organization design.

Our results inform not only the optimal evaluation of employee performance, but are also suggestive about the selection of information sources. Monitoring effort directly remains desirable. Among measures that combine information about effort and ability, the principal prefers measure that are less sensitive to ability. This is in sharp contrast with models of implicit incentives. There, the fact that a signal combines effort and ability is the source of incentives, as the agent exerts effort to avoid being perceived as low-ability. The analysis of evaluation design when both explicit and implicit incentives are present – including the distinction between internal evaluation and externally visible evaluation – is an interesting avenue for future research.

## A Proofs

Some of the proofs allow for a general strictly concave utility function  $u$  with strictly convex inverse  $w$ .

*Proof of Proposition 1:* To see that the mapping from posteriors to wages is 1:1 and deterministic, let  $m(s) := P_\mu^1 p(s|y_H) + (1 - P_\mu^1) p(s|y_L)$  denote the probability of the signal under high effort. It is easy to see that the contracting problem (3) is equivalent to the utility space problem

$$\begin{aligned} \max_{S, p_H, p_L, v} \quad & P_\mu Y + \int_S \left( \delta \Pi_{t+1}(\hat{\mu}(s)) - w(v(s)) \right) m(s) \, ds \\ \text{s.t.} \quad & \int_S v(s) m(s) \, ds - c \geq U \\ & \int_S \left( b + \Delta b \mu \right) \frac{\mu(s) - \mu}{\mu(1 - \mu)(\Delta a + \Delta b)} v(s) m(s) \, ds \geq c \quad (\text{IC}) \\ & \int_S p(s|y_H) \, ds = \int_S p(s|y_L) \, ds = 1 \quad (\text{S}) \end{aligned}$$

where we used the representation of the IC in (6) and the fact that

$$\mu(s) - \mu = \mu(1 - \mu)(\Delta a + \Delta b)(p(s|y_H) - p(s|y_L)).$$

Suppose there are two signals  $s, s'$  with  $\mu(s) = \mu(s')$  and different utilities  $v(s) \neq v(s')$ . We could then set  $\tilde{v} = \frac{m(s)}{m(s)+m(s')} v(s) + \frac{m(s')}{m(s)+m(s')} v(s')$  after both signals. This modification leaves all constraints unchanged, but reduces the costs of incentives since  $w$  is strictly convex.

Therefore, the payoff of any contract is pinned down uniquely by its induced distribution over posterior beliefs and mapping from posteriors to utilities, where optimality allows us to restrict attention to deterministic mappings by the above.

To see the bounds on posteriors, consider

$$\mu(s) = \mu \frac{1 + (a + \Delta a + b + \Delta b) \left( \frac{p(s|y_H)}{p(s|y_L)} - 1 \right)}{1 + (a + \Delta a \mu + b + \Delta b \mu) \left( \frac{p(s|y_H)}{p(s|y_L)} - 1 \right)}$$

This expression is maximized for  $p_L = 0$ , which induces the upper bound, and minimized for  $p_H = 0$ , which induces the lower bound. □

*Proof of Proposition 2:* Note that full information is strictly Blackwell more informative than any other information structure. Then, the result follows from Proposition 13 in [Grossman and Hart \(1983\)](#). Since both the Blackwell comparison as well as the concavity

of the utility function are strict, uniqueness follows from an immediate generalization of their proof.  $\square$

*Proof of Proposition 3:* By standard arguments, both the participation and the incentive constraint are binding. Hence

$$\Pi_2(\mu) = P_\mu Y - P_\mu w\left(U - c + (1 - P_\mu)\frac{c}{b + \Delta b\mu}\right) - (1 - P_\mu)w\left(U - c - P_\mu\frac{c}{b + \Delta b\mu}\right).$$

Note that we require  $U - P_\mu\frac{c}{b + \Delta b\mu} > 0$  to satisfy the implicit nonnegativity constraint in the agent's utility function. Since  $\frac{\partial}{\partial \mu} P_\mu\frac{c}{b + \Delta b\mu} \propto \Delta ab - \Delta ba < 0$ , this is implied by  $U > \frac{a+b}{b}c$ . It is easy to verify that

$$\begin{aligned} \Pi_2''(\mu) &\propto 2(b\Delta a - a\Delta b)(b\Delta a + \Delta b(1 - a))(b + \Delta b\mu) \\ &\quad \cdot \left[ w'\left(U + (1 - P_\mu)\frac{c}{b + \Delta b\mu}\right) - w'\left(U - P_\mu\frac{c}{b + \Delta b\mu}\right) \right] \\ &\quad - cP_\mu(b\Delta a + \Delta b(1 - a))^2 w''\left(U + (1 - P_\mu)\frac{c}{b + \Delta b\mu}\right) \\ &\quad - c(1 - P_\mu)(b\Delta a - a\Delta b)^2 w''\left(U - P_\mu\frac{c}{b + \Delta b\mu}\right) \end{aligned}$$

The two latter terms are clearly negative, and so is the first, since  $b\Delta a - a\Delta b < 0$ . The statement about the Blackwell comparison is immediate.  $\square$

We will prove the Lemmas used in the text to establish Theorem 1 first.

*Proof of Lemma 1:* Note that the optimal wage in the dual problem is given by

$$w^*(\lambda, \hat{\mu}) = \max\left\{0, u'^{-1}\left(\lambda_P + \frac{b + \Delta b\mu}{(\Delta a + \Delta b)\mu(1 - \mu)}(\hat{\mu} - \mu)\right)\right\}$$

Let  $\lambda^*$  denote a pair of multipliers such that the constraints are binding or the respective Lagrange multiplier is zero and the associated wage function is feasible. Then,

$$\begin{aligned} \inf_{\lambda} \sup_w \mathcal{L}(w, \lambda) &= \inf_{\lambda} \mathcal{L}(w^*(\lambda, \cdot), \lambda) \leq \mathcal{L}(w^*(\lambda^*, \cdot), \lambda^*) \\ &= P_\mu^1 Y - \int w^*(\lambda, \hat{\mu}) m(\hat{\mu}) d\hat{\mu} \leq \sup_w \inf_{\lambda} \mathcal{L}(w, \lambda) \end{aligned}$$

whence the problem satisfies strong duality.

To find such a  $\lambda^*$  consider the dual problem. Note that by an envelope argument

$$\frac{\partial \mathcal{L}(w^*(\lambda, \cdot), \lambda)}{\partial \lambda_P} = \int (u(w^*(\lambda, \hat{\mu})) - U - c) m(\hat{\mu}) d\hat{\mu}$$

and similarly for  $\lambda_{IC}$  and hence the Hessian of the objective is given by

$$\begin{pmatrix} \int f(\hat{\mu}) d\hat{\mu} & \int f(\hat{\mu}) \frac{b+\Delta b\mu}{(\Delta a+\Delta b)\mu(1-\mu)} (\hat{\mu}-\mu) d\hat{\mu} \\ \int f(\hat{\mu}) \frac{b+\Delta b\mu}{(\Delta a+\Delta b)\mu(1-\mu)} (\hat{\mu}-\mu) d\hat{\mu} & \int f(\hat{\mu}) \left[ \frac{b+\Delta b\mu}{(\Delta a+\Delta b)\mu(1-\mu)} (\hat{\mu}-\mu) \right]^2 d\hat{\mu} \end{pmatrix} \quad (38)$$

where  $f(\hat{\mu}) := \rho'(\lambda_P + \lambda_{IC} \frac{b+\Delta b\mu}{(\Delta a+\Delta b)\mu(1-\mu)} (\hat{\mu}-\mu))m(\hat{\mu})$  is a positive kernel and the range of integration is over  $\hat{\mu}$  such that  $u'^{-1}(\lambda_P + \frac{b+\Delta b\mu}{(\Delta a+\Delta b)\mu(1-\mu)} (\hat{\mu}-\mu)) \geq 0$ . Hence, the integral  $\int f(\hat{\mu})g_1(\hat{\mu})g_2(\hat{\mu})d\hat{\mu}$  is an inner product (between functions that share support with  $m$ ), the objective is weakly convex by Cauchy-Schwarz, as the determinant of the Hessian reads

$$\langle g_1, g_1 \rangle \langle g_2, g_2 \rangle - \langle g_1, g_2 \rangle^2 \geq 0$$

for  $g_1 = 1$  and  $g_2 = \frac{b+\Delta b\mu}{(\Delta a+\Delta b)\mu(1-\mu)} (\hat{\mu}-\mu)$ .

If  $m$  is nondegenerate, we can bound the space of multipliers in the dual problem without loss. To see this, note that the infimum is bounded above by  $P_\mu^1 Y$  and this bound is achieved by  $\lambda = 0$ . For the first bound, suppose the wage is constant. Then, for an optimal  $\lambda$ , we require that even for this suboptimal  $w$ :

$$P_\mu^1 Y - w + \lambda_P (u(w) - U - c) - \lambda_{IC} c \leq \mathcal{L}(w^*(\lambda, \cdot), \lambda) \leq P_\mu^1 Y$$

Choose  $w$  such that  $u(w) - U - c \geq 2c$ . Then, the above implies that for a suitable  $A$ ,

$$\lambda_P \leq A + \frac{1}{2}\lambda_{IC}$$

Similarly, construct a wage function such that  $\int u(w(\hat{\mu})) \frac{b+\Delta b\mu}{(\Delta a+\Delta b)\mu(1-\mu)} (\hat{\mu}-\mu)m(\hat{\mu}) d\hat{\mu} - c \geq 2(U+c)$ . Then, optimality of  $\lambda$  requires

$$\begin{aligned} & P_\mu^1 Y - w + \lambda_P \left( \int u(w)m(\hat{\mu}) d\hat{\mu} - U - c \right) \\ & + \left( \lambda_{IC} \int u(w) \frac{b+\Delta b\mu}{(\Delta a+\Delta b)\mu(1-\mu)} (\hat{\mu}-\mu)m(\hat{\mu}) d\hat{\mu} - c \right) \leq \\ & P_\mu^1 Y - w - \lambda_P (U+c) + 2(U+c)\lambda_{IC} \leq P_\mu^1 Y \end{aligned}$$

and therefore, for suitable  $B$ ,

$$\lambda_{IC} \leq B + \frac{1}{2}\lambda_P.$$

Taking both inequalities together, it is easy to see that we can bound  $\lambda$  and assume the set of  $\lambda$  to be compact. Hence, the dual problem has a solution. By convexity, it satisfies the FOC or comparative slackness. As the FOC is equal to the respective constraints, the wage at the solution either satisfies the constraints with equality or the constraint is slack and the multiplier is zero. We have found the desired  $\lambda^*$  for the case of nondegenerate  $m$ .

For degenerate  $m$ , the problem is infeasible and hence both the primal and dual value are  $-\infty$ . Therefore the wage problem satisfies strong duality.  $\square$

*Proof of Lemma 2:* Clearly, the space of posterior distributions satisfying (BP) is compact in the weak topology, and, as  $\ell^*$  is continuous and bounded for any  $\lambda$ , the problem is continuous and linear in  $m$ . Continuity in  $\lambda$  is immediate. To see quasi-convexity, note that by an envelope argument

$$\frac{\partial \ell^*}{\partial \lambda_P} = \int \left( \rho(\lambda_P + \lambda_{IC} \frac{b + \Delta b \mu}{(\Delta a + \Delta b) \mu (1 - \mu)}) (\hat{\mu} - \mu) - U - c \right) m(\hat{\mu}) d\hat{\mu}$$

and similarly for  $\lambda_{IC}$  and hence the Hessian of the objective is given by (38), the objective is weakly convex by Cauchy-Schwarz. Therefore, the problem satisfies the conditions of Sion's Minimax Theorem and we have

$$\inf_{\lambda \geq 0} \sup_{w, m \text{ s.t. (BP)}} \mathcal{L}(m, w; (\lambda_P, \lambda_{IC})) = \sup_{m \text{ s.t. (BP)}} \inf_{\lambda \geq 0} \sup_w \mathcal{L}(m, w; (\lambda_P, \lambda_{IC})).$$

$\square$

*Proof of Lemma 3:* From (20), it is easy to see that

$$\frac{\partial^3}{\partial \hat{\mu}^3} \ell^*(\hat{\mu}; \lambda) = \delta \Pi_2'''(\hat{\mu}) + \lambda_{IC}^3 \left[ \frac{b + \Delta b \mu}{(\Delta a + \Delta b) \mu (1 - \mu)} \right]^3 \rho''(\lambda_P + \lambda_{IC} \frac{b + \Delta b \mu}{(\Delta a + \Delta b) \mu (1 - \mu)}) (\hat{\mu} - \mu)$$

It is elementary but tedious to show that

$$\begin{aligned} \Pi_2'''(\mu) = & \frac{c}{(b + \mu \Delta b)^6} \left[ 6 \Delta b (b \Delta a + (1 - a) \Delta b) (a \Delta b - b \Delta a) (b + \mu \Delta b)^2 (w'(u_H) - w'(u_L)) \right. \\ & + 3c (a \Delta b - b \Delta a)^2 (b + \mu \Delta b) (b \Delta a + \Delta b (2 - 2a - b - \mu(\Delta a + \Delta b))) w''(u_L) \\ & + 3c (b \Delta a + (1 - a) \Delta b)^2 (b + \mu \Delta b) (a \Delta b - b \Delta a + \Delta b (a + b + \mu(\Delta a + \Delta b))) w''(u_H) \\ & - c^2 (a \Delta b - b \Delta a)^3 (1 - a - b - \mu(\Delta a + \Delta b)) w'''(u_L) \\ & \left. + c^2 (b \Delta a + (1 - a) \Delta b)^3 (a + b + \mu(\Delta a + \Delta b)) w'''(u_H) \right] \end{aligned}$$

where  $u_L = U - \frac{a + \mu \Delta a}{b + \mu \Delta b} c$  and  $u_H = U + \frac{1 - a - \mu \Delta a}{b + \mu \Delta b} c$ . Under Assumption 1.2, we have  $\Pi_2''' > 0$ . Hence, since  $\rho'' \geq 0$ , we have  $\frac{\partial^3}{\partial \hat{\mu}^3} \ell^*(\hat{\mu}; \lambda) \geq 0$  for all  $\lambda$ .

Let  $\text{cav} f = \max_{\psi, \psi' \in [\underline{\mu}, \bar{\mu}], \alpha \in [0, 1] \text{ s.t. } \alpha \psi + (1 - \alpha) \psi' = \mu} \{\alpha f(\psi) + (1 - \alpha) f(\psi')\}$  denote the concavification of function  $f$  on the interval  $[\underline{\mu}, \bar{\mu}]$  and consider the set of beliefs that can



be used to generate the concavification of  $\ell^*$  at the prior belief  $\mu$ ,

$$\Psi(\lambda_P, \lambda_{IC}) := \{\psi \in [\underline{\mu}, \bar{\mu}] \mid \exists \psi' \in [\underline{\mu}, \bar{\mu}], \alpha \in [0, 1] \text{ s.t. } \alpha\psi + (1 - \alpha)\psi' = \mu \text{ and} \quad (39)$$

$$\text{cav}\ell^*(\mu; \lambda_P, \lambda_{IC}) = \alpha\ell^*(\psi; \lambda_P, \lambda_{IC}) + (1 - \alpha)\ell^*(\psi'; \lambda_P, \lambda_{IC})\} \quad (40)$$

We have to show that the set is at most cardinality two and has the described structure. First, consider the case when  $\ell^*$  is globally concave. Then it is strictly concave at  $\mu$  and, clearly,  $\Psi(\lambda_P, \lambda_{IC}) = \{\mu\}$ . If instead  $\ell^*$  is globally convex, then  $\Psi(\lambda_P, \lambda_{IC}) = \{\underline{\mu}, \bar{\mu}\}$ . In all other cases, there exists a  $\psi$  such that  $\ell^*$  is strictly concave for  $\hat{\mu} < \psi$  and strictly convex for  $\hat{\mu} > \psi$ . Then, the concavification of  $\ell^*$  is equivalent to  $\ell^*$  up to a threshold  $\mu^* < \psi$  and linear, generated by  $\mu^*, \bar{\mu}$  afterwards. Hence, either  $\Psi(\lambda_P, \lambda_{IC}) = \{\mu\}$ , or  $\Psi(\lambda_P, \lambda_{IC}) = \{\mu^*, \bar{\mu}\}$ . The remaining statements are immediate from Bayes plausibility,  $M(\mu^*)\mu^* + M(\bar{\mu})\bar{\mu} = \mu$  and  $M(\mu^*) + M(\bar{\mu}) = 1$ .  $\square$

*Proof of Lemma 4:* To show nondegeneracy in the simplified problem, we need to show that the optimal distribution of posteriors is nondegenerate. Then, existence follows from the bounds on dual multipliers in Lemma 1. To this purpose, we show that there exists an  $\bar{\epsilon}_1$  such that  $\mu^* < \mu - \bar{\epsilon}_1$ . This also establishes non-degeneracy of the optimal information structure.

Suppose not, let  $\mu^* = \mu - \epsilon$  and we will show that the costs of providing incentives diverge as  $\epsilon \rightarrow 0$ . To see this, note that

$$\begin{aligned} m(\mu^*)\mu^* + m(\bar{\mu})\bar{\mu} &= \mu \\ m(\bar{\mu}) &= \frac{\mu - m(\mu^*)\mu^*}{\bar{\mu} - \mu} \\ &= \frac{\mu - m(\mu^*)\mu - \epsilon}{\bar{\mu} - \mu} \\ &= \epsilon \frac{m(\mu^*)}{\bar{\mu} - \mu} \leq \epsilon \frac{1}{\bar{\mu} - \mu} \end{aligned}$$

In the IC constraint, we have

$$\begin{aligned} c &\leq \frac{b + \Delta b\mu}{(\Delta a + \Delta b)\mu(1 - \mu)} [m^*(\mu^* - \mu)u(w^*) + (1 - m^*)(\bar{\mu} - \mu)u(\bar{w})] \\ &\leq \frac{b + \Delta b\mu}{(\Delta a + \Delta b)\mu(1 - \mu)} \left[ \epsilon \frac{1}{\bar{\mu} - \mu} (\bar{\mu} - \mu)(u(\bar{w}) - u(\underline{w})) \right] \end{aligned}$$

Hence, as  $\epsilon \rightarrow 0$ , we require  $u(\bar{w}) \geq c_0\epsilon^{-1}$ , for a suitable constant  $c_0$ . But then, the objective is  $\leq c_1 - \epsilon \cdot c_2w(\epsilon^{-1}) \rightarrow -\infty$  for suitable constants, which is clearly not optimal.

Hence, the optimal distribution of posteriors is nondegenerate and a solution exists. Is is unique since the problem is concave with convex constraint sets.  $\square$

*Proof of Theorem 1:* The claims about the information structure follow immediately from the previous lemma.  $\square$

*Proof of Proposition 4:* We rewrite the simplified problem, noting that  $m^* = \frac{\bar{\mu} - \mu}{\bar{\mu} - \mu^*}$  and maximizing out wages. First, note that the IC constraint reads

$$\begin{aligned} & \frac{b + \Delta b\mu}{(\Delta a + \Delta b)\mu(1 - \mu)} [m^*(\mu^* - \mu)u(w^*) + (1 - m^*)(\bar{\mu} - \mu)u(\bar{w})] = \\ & \frac{b + \Delta b\mu}{(\Delta a + \Delta b)\mu(1 - \mu)} \left[ \frac{\bar{\mu} - \mu}{\bar{\mu} - \mu^*} (\mu^* - \mu)u(w^*) + \frac{\mu - \mu^*}{\bar{\mu} - \mu^*} (\bar{\mu} - \mu)u(\bar{w}) \right] = \\ & \frac{b + \Delta b\mu}{\Delta b\mu(1 - \mu)} \frac{(\bar{\mu} - \mu)(\mu - \mu^*)}{\bar{\mu} - \mu^*} [u(\bar{w}) - u(w^*)] \geq c \end{aligned}$$

Then  $u(\bar{w}) = \lambda_P + \lambda_{IC} \frac{b + \Delta b\mu}{(\Delta a + \Delta b)\mu(1 - \mu)} (\bar{\mu} - \mu)$  and  $u(w^*) = \lambda_P - \lambda_{IC} \frac{b + \Delta b\mu}{(\Delta a + \Delta b)\mu(1 - \mu)} (\mu - \mu^*)$ . The multipliers are  $\lambda_P = U + c$  and

$$\lambda_{IC} = \frac{c}{\left( \frac{b + \Delta b\mu}{(\Delta a + \Delta b)\mu(1 - \mu)} \right)^2 (\bar{\mu} - \mu)(\mu - \mu^*)}$$

By an envelope argument, the first order condition for  $\mu^*$  is (writing in utility space)

$$\begin{aligned} 0 &= \delta \left[ \frac{\bar{\mu} - \mu}{(\bar{\mu} - \mu^*)^2} (\Pi_2(\mu^*) - \Pi_2(\bar{\mu})) + \frac{\bar{\mu} - \mu}{\bar{\mu} - \mu^*} \Pi_2'(\mu^*) \right] + \\ & \frac{1}{2} \left[ \frac{\bar{\mu} - \mu}{(\bar{\mu} - \mu^*)^2} (u^{*2} - \bar{u}^2) + \frac{\bar{\mu} - \mu}{\bar{\mu} - \mu^*} \lambda_{IC} \left( \frac{b + \Delta b\mu}{\Delta b\mu(1 - \mu)} \right) u^* \right] = \\ \delta & \left[ \frac{\bar{\mu} - \mu}{(\bar{\mu} - \mu^*)^2} (\Pi_2(\mu^*) - \Pi_2(\bar{\mu})) + \frac{\bar{\mu} - \mu}{\bar{\mu} - \mu^*} \Pi_2'(\mu^*) \right] - \frac{1}{2} \lambda_{IC}^2 \left( \frac{b + \Delta b\mu}{(\Delta a + \Delta b)\mu(1 - \mu)} \right)^2 (\bar{\mu} - \mu) \end{aligned}$$

as is straightforward but tedious to show. Plugging in for the multiplier and multiplying through, we arrive at the condition

$$\delta [\Pi_2(\mu^*) - \Pi_2(\bar{\mu}) + (\bar{\mu} - \mu^*) \Pi_2'(\mu^*)] = \frac{1}{2} \left( \frac{(\Delta a + \Delta b)\mu(1 - \mu)}{b + \Delta b\mu} c \right)^2 \frac{(\bar{\mu} - \mu^*)^2}{(\bar{\mu} - \mu)^2 (\mu - \mu^*)^2}$$

Note that this condition holds for an interior solution. As  $\mu^* \rightarrow \mu$ , the RHS diverges while LHS stays bounded, so there will never be a corner solution at this limit. As  $\mu^* \rightarrow \bar{\mu}$ ,

LHS grows as  $\Pi_2'' < 0$  and RHS shrinks, but both stay bounded. We therefore have a corner solution at  $\mu^* = \underline{\mu}$  if (13) is violated.  $\square$

*Proof of Proposition 6:* To see the first statement, consider the problem in signal/utility space. Then, the cost of incentives is

$$\int ((a + b + \eta(\Delta a + \Delta b))p_H(s) + (1 - a - b - \eta(\Delta a + \Delta b))p_L(s)) w(u(s)) ds$$

and the constraints depend on

$$\int ((a + b + \mu(\Delta a + \Delta b))p_H(s) + (1 - a - b - \mu(\Delta a + \Delta b))p_L(s)) u(s) ds$$

and similar for IC. For a given  $p_L, p_S, u$ , we will construct a cheaper fully informative contract. Consider providing  $\int p_H(s)u(s) ds$  for certain after high output and  $\int p_L(s)u(s) ds$  after low output. The constraints are unchanged, so this contract is feasible. It is also cheaper by the convexity of  $w$ , strictly so if  $p_L, p_H$  were not degenerate.

To see monotonicity and concavity, note that

$$\begin{aligned} \Pi_2^\eta(\mu) &= (a + b + \eta(\Delta a + \Delta b))Y - (a + b + \eta(\Delta a + \Delta b))\frac{1}{2} \left( U + (1 - P_\mu)\frac{c}{b + \Delta b\mu} \right)^2 \\ &\quad - (1 - a - b - \eta(\Delta a + \Delta b))\frac{1}{2} \left( U - P_\mu\frac{c}{b + \Delta b\mu} \right)^2 \\ \Pi_2^{0'}(\mu) &= \frac{c}{(b + \Delta b\mu)^3} [ac\Delta b(1 - a - b) + bc\Delta b(1 - a - b - \mu(\Delta a + \Delta b)) \\ &\quad - bc\Delta a(\Delta a + \Delta b)\mu + b(\Delta a + \Delta b)(b + \Delta b\mu)U] \\ &> \frac{c}{(b + \Delta b\mu)^3} [c\Delta b(1 - a - b - \mu(\Delta a + \Delta b))(a + b) + b(\Delta a + \Delta b)(b + \Delta b\mu)U] \\ &> 0 \\ \Pi_2^{0''}(\mu) &= -\frac{c}{(b + \Delta b\mu)^4} [c\Delta b^2(3a(1 - a - b) + b(3 - 3a - 2b - 2\Delta b\mu)) \\ &\quad + 2b\Delta b(\Delta a + \Delta b)(b + \Delta b\mu)U - cb(\Delta a^2 + 2\Delta a\Delta b)(2\Delta b\mu - b)] \\ &< -\frac{c}{(b + \Delta b\mu)^4} [c\Delta b^2(3a(1 - a - b) + b(3 - 3a - 2b - 2\Delta b\mu)) \\ &\quad + 2b\Delta b(\Delta a + \Delta b)(b + \Delta b\mu)U - cb(\Delta a^2 + 2\Delta a\Delta b)(2\Delta b\mu - b)] \\ &< 0 \end{aligned}$$

using the fact that either  $(2\Delta b\mu - b)$  is negative or we can use log-supermodularity. The results for overconfidence also follow from straightforward but tedious computation.

To see the result on information, note that information corresponds to a MPS of  $m$ , which the principal evaluates as an integral of  $\left[\eta \frac{\hat{\mu}}{\mu} + (1 - \eta) \frac{1 - \hat{\mu}}{1 - \mu}\right] \Pi_2^\eta(\hat{\mu})$ .

If  $\eta = 0$ : Since  $U > \frac{a+b}{b}c$ ,  $\frac{\partial^2}{\partial \hat{\mu}^2} \left(\left[\eta \frac{\hat{\mu}}{\mu} + (1 - \eta) \frac{1 - \hat{\mu}}{1 - \mu}\right] \Pi_2^\eta(\hat{\mu})\right) > 0$  follows from direct computation.

If  $\Delta b < 0$ : The sign is ambiguous,

$$\begin{aligned} \frac{\partial^2}{\partial \hat{\mu}^2} \left( \left[ \eta \frac{\hat{\mu}}{\mu} + (1 - \eta) \frac{1 - \hat{\mu}}{1 - \mu} \right] \Pi_2^\eta(\hat{\mu}) \right) &\propto c \left[ b^2 \Delta a (2b - \Delta b(6 + 4\mu) - \Delta a(6 + 3\mu)) \right. \\ &\quad + (2b - \Delta b\mu) (\Delta b(1 - a)(b + \Delta a + \Delta b) + 4b\Delta a(\Delta a + \Delta b) \\ &\quad + \Delta ba(1 - a - b - \Delta a - \Delta b))] \\ &\quad \left. + 2b(b + \Delta b)(\Delta a + \Delta b)(b + \Delta b\mu)(U - c) \right] \end{aligned}$$

In particular, the expression is increasing in  $U$ . Furthermore, let  $U = \frac{a+b}{b}c - \delta$  in order to ensure that the constraint is satisfied as we change  $b$ . For  $b = 0$ , we get

$$\frac{\partial^2}{\partial \hat{\mu}^2} \left( \left[ \eta \frac{\hat{\mu}}{\mu} + (1 - \eta) \frac{1 - \hat{\mu}}{1 - \mu} \right] \Pi_2^\eta(\hat{\mu}) \right) \propto -[a(1 - a) + (\Delta a + \Delta b)(1 - 4a)]$$

and if this expression is negative, we get the cutoff by continuity. □

*Proof of Theorem 2:* As this proof closely follows the same template as the proof of Theorem 1, we will be brief. All functions relate to Section 4, we refrain from using decorators to mark this association.

*(Optimal Wages)* The pointwise optimal wage schedule in the Lagrangian associated with (29) is

$$w^*(\hat{\mu}, \lambda) = \frac{1}{2} \left( \frac{1}{\eta \frac{\hat{\mu}}{\mu} + (1 - \eta) \frac{1 - \hat{\mu}}{1 - \mu}} \right)^2 \left( \lambda_P + \lambda_{IC} \frac{b + \Delta b\mu}{(\Delta a + \Delta b)\mu(1 - \mu)} (\hat{\mu} - \mu) \right)^2$$

*(Info Design)* The Lagrangian is additively separable and

$$\begin{aligned} \ell^*(\hat{\mu}; \lambda) = &P_0 Y + \left[ \eta \frac{\hat{\mu}}{\mu} + (1 - \eta) \frac{1 - \hat{\mu}}{1 - \mu} \right] [\delta \Pi_2(\hat{\mu}) - w^*(\hat{\mu}, \lambda)] + \lambda_P (u(w^*(\hat{\mu}, \lambda)) - c - U) \\ &+ \lambda_{IC} \left( \frac{b + \Delta b\mu}{(\Delta a + \Delta b)\mu(1 - \mu)} (\hat{\mu} - \mu) u(w^*(\hat{\mu}, \lambda)) - c \right) \end{aligned}$$

If  $\eta = 0$ : Then,

$$\begin{aligned}\frac{\partial^2}{\partial \hat{\mu}^2} \ell^*(\hat{\mu}; \lambda) &= \frac{1 - \mu}{(1 - \hat{\mu})^3} \left( \lambda_P + \lambda_{IC} \frac{b + \Delta b \mu}{(\Delta a + \Delta b) \mu} \right)^2 + \delta \frac{(1 - \hat{\mu}) \Pi_2^{0''}(\hat{\mu}) - 2 \Pi_2^{0'}(\hat{\mu})}{1 - \mu} \\ \frac{\partial^3}{\partial \hat{\mu}^3} \ell^*(\hat{\mu}; \lambda) &= 3 \frac{1 - \mu}{(1 - \hat{\mu})^4} \left( \lambda_P + \lambda_{IC} \frac{b + \Delta b \mu}{(\Delta a + \Delta b) \mu} \right)^2 + \delta \frac{(1 - \hat{\mu}) \Pi_2^{0''' }(\hat{\mu}) - 3 \Pi_2^{0''}(\hat{\mu})}{1 - \mu}\end{aligned}$$

and  $\frac{\partial^3}{\partial \hat{\mu}^3} \ell^*(\hat{\mu}; \lambda) > 0$ . Lemma 3 goes through. We can apply the proof of Lemmas 4 and 2 mutatis mutandis and arrive at the Theorem.

If  $\eta = 1$ : Then, we have

$$\begin{aligned}\frac{\partial^2}{\partial \hat{\mu}^2} \ell^*(\hat{\mu}; \lambda) &= \frac{\mu}{\hat{\mu}^3} \left( \lambda_P - \lambda_{IC} \frac{b + \Delta b \mu}{(\Delta a + \Delta b)(1 - \mu)} \right)^2 + \delta \frac{\hat{\mu} \Pi_2^{1''}(\hat{\mu}) + 2 \Pi_2^{1'}(\hat{\mu})}{\mu} \\ \frac{\partial^3}{\partial \hat{\mu}^3} \ell^*(\hat{\mu}; \lambda) &= -3 \frac{\mu}{\hat{\mu}^4} \left( \lambda_P - \lambda_{IC} \frac{b + \Delta b \mu}{(\Delta a + \Delta b)(1 - \mu)} \right)^2 + \delta \frac{\hat{\mu} \Pi_2^{\eta''' }(\hat{\mu}) + 3 \Pi_2^{1''}(\hat{\mu})}{\mu}\end{aligned}$$

It is straightforward but tedious to show that  $\frac{\partial^2}{\partial \hat{\mu}^2} \ell^*(\hat{\mu}; \lambda) = 0 \implies \frac{\partial^3}{\partial \hat{\mu}^3} \ell^*(\hat{\mu}; \lambda) < 0$  and therefore the Lagrangian is either convex or convex to concave (it cannot be globally concave by incentive compatibility). Hence, a lenient information structure is optimal and the lemmas generalize. □

## References

- Adrian, Tobias and Mark M. Westerfield**, “Disagreement and Learning in a Dynamic Contracting Model,” *Review of Financial Studies*, October 2009, 22 (10), 3873–3906. 7
- Alonso, Ricardo and Odilon Câmara**, “Bayesian Persuasion with Heterogeneous Priors,” *Journal of Economic Theory*, September 2016, 165, 672–706. 8, 26
- Aumann, Robert J. and Michael Maschler**, *Repeated Games with Incomplete Information*, Cambridge, Mass: MIT Press, 1995. 18
- Bénabou, Roland and Jean Tirole**, “Belief in a Just World and Redistributive Politics,” *The Quarterly Journal of Economics*, May 2006, 121 (2), 699–746. 24
- Bergemann, Dirk and Stephen Morris**, “Information Design: A Unified Perspective,” *Journal of Economic Literature*, March 2019, 57 (1), 44–95. 8

- Bhaskar, V. and George J. Mailath**, “The Curse of Long Horizons,” *Journal of Mathematical Economics*, May 2019, *82*, 74–89. 7, 33
- Boleslavsky, Raphael and Kyungmin Kim**, “Bayesian Persuasion and Moral Hazard,” *SSRN Electronic Journal*, 2017. 8, 13, 18
- Boretz, Elizabeth**, “Grade Inflation and the Myth of Student Consumerism,” *College Teaching*, 2004, *52* (2), 42–46. 25
- Brunnermeier, Markus K. and Martin Oehmke**, “The Maturity Rat Race,” *The Journal of Finance*, 2013, *68* (2), 483–521. 31
- Burks, Stephen V., Jeffrey P. Carpenter, Lorenz Goette, and Aldo Rustichini**, “Overconfidence and Social Signalling,” *The Review of Economic Studies*, 2013, *80* (3 (284)), 949–983. 24
- Charness, Gary and Uri Gneezy**, “Portfolio Choice and Risk Attitudes: An Experiment,” *Economic Inquiry*, 2010, *48* (1), 133–146. 24
- Cho, In-Koo and David M. Kreps**, “Signaling Games and Stable Equilibria,” *The Quarterly Journal of Economics*, May 1987, *102* (2), 179–221. 32
- **and Joel Sobel**, “Strategic Stability and Uniqueness in Signaling Games,” *Journal of Economic Theory*, April 1990, *50* (2), 381–413. 32
- Datar, Srikant, Susan Cohen Kulp, and Richard A. Lambert**, “Balancing Performance Measures,” *Journal of Accounting Research*, June 2001, *39* (1), 75–92. 6
- de la Rosa, Leonidas Enrique**, “Overconfidence and Moral Hazard,” *Games and Economic Behavior*, November 2011, *73* (2), 429–451. 7
- Demarzo, Peter M. and Yuliy Sannikov**, “Learning, Termination, and Payout Policy in Dynamic Incentive Contracts,” *The Review of Economic Studies*, January 2017, *84* (1), 182–236. 7, 33
- Dewatripont, Mathias, Ian Jewitt, and Jean Tirole**, “The Economics of Career Concerns, Part I: Comparing Information Structures,” *The Review of Economic Studies*, January 1999, *66* (1), 183–198. 7

- Dittmann, Ingolf and Ernst Maug**, “Lower Salaries and No Options? On the Optimal Structure of Executive Pay,” *The Journal of Finance*, 2007, *62* (1), 303–343. 9
- Doval, Laura and Vasiliki Skreta**, “Constrained Information Design: Toolkit,” Technical Report 2018. 8, 16
- Dye, Ronald A.**, “Optimal Monitoring Policies in Agencies,” *The RAND Journal of Economics*, 1986, *17* (3), 339–350. 6
- Ederer, Florian**, “Feedback and Motivation in Dynamic Tournaments,” *Journal of Economics & Management Strategy*, September 2010, *19* (3), 733–769. 5
- , **Richard Holden, and Margaret Meyer**, “Gaming and Strategic Opacity in Incentive Provision,” *The RAND Journal of Economics*, December 2018, *49* (4), 819–854. 5
- Edmans, Alex and Xavier Gabaix**, “The Effect of Risk on the CEO Market,” *The Review of Financial Studies*, August 2011, *24* (8), 2822–2863. 9
- Ekmekci, Mehmet and Nenad Kos**, “Signaling Covertly Acquired Information,” Technical Report 2019. 32
- Ely, Jeffrey and Martin Szydlowski**, “Moving the Goalposts,” *Journal of Political Economy*, May 2019, p. 704387. 6
- Fang, Hanming and Giuseppe Moscarini**, “Morale Hazard,” *Journal of Monetary Economics*, May 2005, *52* (4), 749–777. 6, 8
- Feltham, Gerald A. and Jim Xie**, “Performance Measure Congruity and Diversity in Multi-Task Principal/Agent Relations,” *The Accounting Review*, 1994, *69* (3), 429–453. 6
- Filippin, Antonio and Paolo Crosetto**, “Click’n’Roll: No Evidence of Illusion of Control,” *De Economist*, September 2016, *164* (3), 281–295. 24
- Fuchs, William**, “Contracting with Repeated Moral Hazard and Private Evaluations,” *American Economic Review*, September 2007, *97* (4), 1432–1448. 5

- Galperti, Simone**, “Persuasion: The Art of Changing Worldviews,” *American Economic Review*, March 2019, 109 (3), 996–1031. 26
- Georgiadis, George and Balazs Szentes**, “Optimal Monitoring Design,” *Econometrica*, forthcoming, p. 55. 6, 8, 17, 18
- Giat, Yahel, Steve T. Hackman, and Ajay Subramanian**, “Investment under Uncertainty, Heterogeneous Beliefs, and Agency Conflicts,” *Review of Financial Studies*, April 2010, 23 (4), 1360–1404. 7
- Gittleman, Maury and Brooks Pierce**, “How Prevalent Is Performance-Related Pay in the United States? Current Incidence and Recent Trends,” *National Institute Economic Review*, November 2013, 226 (1), R4–R16. 2
- Gonzalez-Hernandez, G., A. Sarker, K. O’Connor, and G. Savova**, “Capturing the Patient’s Perspective: A Review of Advances in Natural Language Processing of Health-Related Text,” *Yearbook of Medical Informatics*, August 2017, 26 (1), 214–227. 2
- Grossman, Sanford J. and Oliver D. Hart**, “An Analysis of the Principal-Agent Problem,” *Econometrica*, January 1983, 51 (1), 7. 2, 5, 14, 37
- Hart, Oliver and Jean Tirole**, “Vertical Integration and Market Foreclosure,” *Brookings Papers on Economic Activity. Microeconomics*, 1990, 1990, 205. 31
- Hoffmann, Florian, Roman Inderst, and Marcus M. Opp**, “Only Time Will Tell: A Theory of Deferred Compensation,” Technical Report 2019. 6
- Holmström, Bengt**, “Moral Hazard and Observability,” *The Bell Journal of Economics*, April 1979, 10 (1), 74–91. 2, 5
- , “Managerial Incentive Problems: A Dynamic Perspective,” *The Review of Economic Studies*, January 1999, 66 (1), 169–182. 2, 6
- **and Paul Milgrom**, “Multitask Principal-Agent Analyses: Incentive Contracts, Asset Ownership, and Job Design,” *Journal of Law, Economics, & Organization*, January 1991, 7, 24–52. 2, 5
- Hörner, Johannes and Nicolas S. Lambert**, “Motivational Ratings,” *Review of Economic Studies*, forthcoming. 2, 7



- Huffman, D., Collin Raymond, and Julia Shvets**, “Persistent Overconfidence and Biased Memory: Evidence from Managers,” Technical Report 2019. 24
- Hügelschäfer, Sabine and Anja Achtziger**, “On Confident Men and Rational Women: It’s All on Your Mind(Set),” *Journal of Economic Psychology*, April 2014, 41, 31–44. 24
- Jehiel, Philippe**, “On Transparency in Organizations,” *The Review of Economic Studies*, April 2015, 82 (2), 736–761. 5
- Kamenica, Emir and Matthew Gentzkow**, “Bayesian Persuasion,” *American Economic Review*, October 2011, 101 (6), 2590–2615. 8, 18
- Kim, Son Ku**, “Efficiency of an Information System in an Agency Model,” *Econometrica*, 1995, 63 (1), 89–102. 5
- Langer, Ellen J.**, “The Illusion of Control,” *Journal of Personality and Social Psychology*, 1975, 32 (2), 311–328. 24
- Larwood, Laurie and William Whittaker**, “Managerial Myopia: Self-Serving Biases in Organizational Planning,” *Journal of Applied Psychology*, April 1977, 62 (2), 194–198. 24
- Lerner, Melvin J.**, *The Belief in a Just World: A Fundamental Delusion*, Boston, MA: Springer US : Imprint : Springer, 1980. 24
- Li, Anqi and Ming Yang**, “Optimal Incentive Contract with Endogenous Monitoring Technology,” *Theoretical Economics*, forthcoming. 6
- Lizzeri, Alessandro, Margaret A. Meyer, and Nicola Persico**, “The Incentive Effects of Interim Performance Evaluations,” Technical Report, Penn Economics Department September 2002. 5
- MacLeod, W. Bentley**, “Optimal Contracting with Subjective Evaluation,” *The American Economic Review*, 2003, 93 (1), 216–240. 5
- Nafziger, Julia**, “Timing of Information in Agency Problems with Hidden Actions,” *Journal of Mathematical Economics*, December 2009, 45 (11), 751–766. 5

- Niederle, Muriel and Lise Vesterlund**, “Do Women Shy Away From Competition? Do Men Compete Too Much?,” *The Quarterly Journal of Economics*, August 2007, 122 (3), 1067–1101. 24
- Orlov, Dmitry, Andrzej Skrzypacz, and Pavel Zryumov**, “Persuading the Principal To Wait,” *Journal of Political Economy*, forthcoming. 31
- Park, Young Joon and Luís Santos-Pinto**, “Overconfidence in Tournaments: Evidence from the Field,” *Theory and Decision*, July 2010, 69 (1), 143–166. 24
- Prat, Julien and Boyan Jovanovic**, “Dynamic Contracts When the Agent’s Quality Is Unknown: Dynamic Contracts,” *Theoretical Economics*, September 2014, 9 (3), 865–914. 7, 33
- Singer, Natasha**, “In a Mood? Call Center Agents Can Tell,” *The New York Times*, October 2013. 2
- Smolin, Alex**, “Dynamic Evaluation Design,” Technical Report 2017. 6
- Treust, Maël Le and Tristan Tomala**, “Persuasion with Limited Communication Capacity,” *Journal of Economic Theory*, November 2019, 184, 104940. 8, 16, 18
- Yaouanq, Yves Le and Peter Schwardmann**, “Learning about One’s Self,” Technical Report 2019. 4