

Discussion Paper Series – CRC TR 224

Discussion Paper No. 358  
Project B 02

Vague by Design:  
Performance Evaluation and Learning From Wages

Franz Ostrizek<sup>1</sup>

June 2022

<sup>1</sup> University of Bonn, Email: [franz.ostrizek@gmail.com](mailto:franz.ostrizek@gmail.com)

Funding by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation)  
through CRC TR 224 is gratefully acknowledged.

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

Franz Ostrizek<sup>†</sup>

May 29, 2022

## 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 is overconfident. Then, 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 to obtain. For example, the availability of board computers and GPS tracking allows for better monitoring

---

\*I am grateful to Roland Bénabou, Pietro Ortoleva, Wolfgang Pesendorfer and Leeat Yariv for their guidance and to Ludmila Matysková, Sofia Moroni, Leon Musolf, Pellumb Reshidi, Evgenii Safonov, Denis Shishkin, Nikhil Vellodi, Can Urgan and seminar audiences at Princeton University, PSE, Northwestern Kellogg, Notre Dame, University of Pittsburgh, ESMT, DICE, University of Bonn, University of Vienna, CERGE-EI, University of Pennsylvania, Sciences Po, the World Congress of the Econometric Society 2020, the World Congress of Game Theory 2020, and the 2021 SEA Annual Meeting for helpful comments and discussions. Funding by the Deutsche Forschungsgemeinschaft (DFG) through CRC TR 224 (Project B02) is gratefully acknowledged.

<sup>†</sup>Department of Economics, University of Bonn; briq – Institute on Behavior and Inequality. franz.ostrizek@gmail.com

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

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? What should be measured and how should the information be aggregated? The theory of incentives seems to offer a simple answer to these questions. Providing incentives to workers is costly exactly 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 not 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 reason is that performance evaluation provides not only the basis of incentives but also shapes learning. By basing wages on a precise evaluation of output, the firm reveals information about output and hence a worker's ability. To illustrate, suppose a worker's contract promises a bonus upon a sufficiently high customer satisfaction score. The worker may gain little information about customer satisfaction through his job directly. With such an evaluation scheme, however, he will know that he cleared the threshold if he receives the bonus, while he learns that he fell short of it if he does not. When the firm prefers workers to remain uninformed, it cannot base incentive pay on the agent's performance directly but needs to use a noisy signal instead.

Why would it be profitable to conceal information about their performance from workers? The costs of providing incentives itself provides a rationale, since uninformed workers are on average cheaper to motivate. When effort and the worker's match-specific ability 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 impact 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>

---

<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-Gaussian career concerns models Holmström (1999); Hörner and Lambert (2021). We discuss some settings in which the addition of noise can be beneficial in the related literature.

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

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 average 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 trade-off. 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 complexity 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, he may receive a low evaluation and thus fail to obtain the bonus. After low quality output, he never obtains the bonus. This information structure is optimal because it 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 effects we study highlight an important consideration in the design of performance evaluation: it is not only the basis for incentives but also for workers’ learning about their ability, their task, and their match to the organization. This is a rich interaction with many facets. Our model is intentionally stylized to retain tractability even though evaluations and contracts are flexibly designed and focuses on learning about match-specific ability. Nevertheless, our results can speak to several notable patterns in incentive systems. First, we would expect tough evaluations in professions with a strong complementarity between skills and effort such as law or investment banking. Indeed, these professions are generally associated with an uncompromising mentality. This is especially prevalent in evaluating the work of fresh associates, which is in line with our model as well: as the tradeoff is intertemporal, evaluations become more informative and less tough over time (Section 5.4). Another feature of evaluations that is in line with our model is that they often focus on the conduct of workers on the job and less so on available measures of output. In our setting, this would be optimal since such evaluations allow the firm to motivate without revealing information about ability.

---

<sup>5</sup>Of course, this mechanism is predicated on two background conditions: First, the agent observes his wages. 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 issues in more detail after introducing the model formally in Section 2 and provide extensions in Section 5.

The influence of explicit incentives on the agent’s learning and confidence becomes even more essential when the agent’s assessments are initially misguided. Indeed, the learning environment then also shapes the evolution of the average belief, 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 (Larwood and Whittaker, 1977; Burks et al., 2013; Huffman et al., 2019), the degree of control they have over their environment (Langer, 1975), or the extent to which they live in a “just world” that rewards effort in the long run (Lerner, 1980). 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 to a continuation value. We also study the long-run evolution of beliefs. Section 6 concludes. The proofs not given in the text are collected in the Appendix.

## Related Literature

This paper contributes to the large literature on information in moral-hazard models. We offer a counterpoint to the classic results establishing that more precise evaluation 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.

We show that noisy evaluation is optimal even though verifiable information about the agent’s true performance would be available. Several strands of the literature show that coarse or noisy evaluation is optimal when such information is not available, for instance 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). Similarly, coarse rewards emerge when the evaluation is subjective, i.e. based on unverifiable private information of the principal (MacLeod, 2003; Fuchs, 2007).

---

<sup>6</sup>We retain the assumption that the agent is Bayesian. In the domain of self-control, there is evidence that individuals update their initially optimistic beliefs in a rational manner Yaouanq and Schwardmann (forthcoming).

That more information about the technology can reduce profits in a moral hazard setting has been noted in the literature in several settings. In general, it is well understood that ex-post incentive compatibility is more demanding than ex-ante incentive compatibility. [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>

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> [Demougin and Fluet \(2001\)](#) study the case with limited liability which results in a binary evaluation. In [Georgiadis and Szentes \(2020\)](#) and [Li and Yang \(2020\)](#) the costs of information acquisition are assumed as part of the technology. [Dai et al. \(2022\)](#) study the optimal contract when the principal can allocate attention between finding good and bad news. [Hoffmann et al. \(2021\)](#) 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. Perhaps the closest reference regarding the analyzed trade-off is [Orlov \(forthcoming\)](#), which studies the optimal contract and intensity of monitoring in a dynamic setting with limited liability. The central trade-off is between monitoring to avoid wasteful investment and thereby revealing information about the continuation value which is costly. We analyze not only the intensity but also the shape of the optimal evaluation when monitoring is required for incentives and its costs stem from the agent’s learning about his type.

Learning about a persistent state and information design are combined in a growing literature. Most closely related to our moral-hazard setting are [Smolin \(2021\)](#) and [Ely and Szydlowski \(2019\)](#) in which the principal uses information, which is valuable for the agent, as an incentive. We analyze the role of information design when – since the principal sets

---

<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.

<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.

incentives and ability is match-specific – 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). For career concerns, it is essential that ability 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. In our setting with explicit incentives, such entangled information is the source of the friction. As a consequence, the role of information is fundamentally different, a point we return to in the conclusion. Hörner and Lambert (2021) analyze the optimal evaluation 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 ability.<sup>10</sup>

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 while we study learning based on an information structure that is designed endogenously by the principal. 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, which is often considerably more tractable (see Bhaskar and Mailath (2019) and Bhaskar (2021) for notable exceptions).

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 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 recent contributions to 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, 2020) – in our case wages. We consider 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 can be analyzed using the transformation approach of Alonso and Câmara (2016).

---

<sup>10</sup>They also show that it 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. Rodina (2018) provides conditions for additional information on output and less prior information about ability to increase effort.

## 2 The Model

A principal (she) employs an agent (he) for two periods. The principal is risk neutral, the agent is risk averse with a strictly increasing utility index  $u : [0, \infty) \rightarrow [0, \infty)$  which we assume to be unbounded.<sup>11</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  $c \cdot e_t$ , 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:

\ 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>12</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  $(S, p, w)$ , comprising an arbitrary measurable signal space  $S$ , a distribution  $p(\cdot|y) \in \Delta(S)$ <sup>13</sup> over signals conditional on high (resp. low) output, and a mapping  $w : S \rightarrow \mathbb{R}$  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.

<sup>11</sup>This is for simplicity to avoid corner solutions.

<sup>12</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.

<sup>13</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 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  $x$ .

<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.

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.

## Discussion

An important feature of our model is that the firm cannot engage in complete backloading of information while still providing incentives. If the principal could record the output of the agent without revealing this information and credibly commit to contingent payments at the end of the relationship and if such a delay was costless (as with risk-neutrality and common discounting), a fully informative and fully delayed evaluation would be optimal.<sup>17</sup> Our mechanism comes into effect when such complete backloading is costly or infeasible. In the main sections, this is achieved in a tractable manner by restricting the principal to period-by-period contracts. This lack of intertemporal commitment ensures that incentives for effort have to be provided in the concurrent period. Wages in the first period are thus informative about output. More generally, risk aversion itself makes backloading information costly: Period-by-period contracting makes the model tractable, but is not the source of the trade-off between incentive provision and learning. We return to this commitment assumption and how it can be relaxed in Section 5.3.

In the design of the evaluation, we allow the principal to reveal more information about output to the agent than is used in wage-setting. This will not be essential, but is allowed for additional generality. What is crucial, however, is that the agent observes at least as much information as is contained in the wages.

The present model features learning based on an endogenous signal distribution with hidden actions and therefore 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 in the baseline model ensure that no endogenous private information arises. Regarding the principal, she does not acquire private information about the type of the agent, since she does not privately observe the quality of output itself but only the result of the public evaluation. We relax this assumption in Section 5.1 and show that our results generalize to the unique

---

<sup>16</sup>In a large organization, this assumption is consistent with the firm observing statistics such as output, profit, and sales in aggregate, as long as it is difficult to link aggregate shortfalls 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.

<sup>17</sup>An instructive comparison is with Orlov (forthcoming): Due to common discounting and risk-neutrality, the cost of incentives is minimized when all payments and also all information about performance is delayed until the end of the relationship. Monitoring instead has an efficiency benefit (weeding out bad projects), which drives the main trade-off. In our setting of risk aversion, by contrast, the cost of incentives is minimized when all information about performance is revealed as early as possible and payments are distributed over time, while fully delayed information revelation would be best to avoid the costs associated with learning about match-specific ability. Period-by-period contracting makes this trade-off more tractable, but is not its source.

equilibrium in a natural class. Regarding the agent, the benchmark model ensures that his posterior belief remains common knowledge after a deviation to lower effort, since effort is noncontractible but observed. If effort was unobserved, the agent would acquire private information about his belief after a deviation and double deviations to low effort in both periods may be 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

In this section, we find the optimal contracts  $(S, p, w)$  by transforming the contracting problem into the space of posterior beliefs and solving it backwards from the second to the first period. We suppress time indices throughout when no confusion can arise.

#### 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 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

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

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

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

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

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

This is a standard moral hazard problem 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*  $\Pi_{t+1}(\mu(s))$ .

**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 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

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

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

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. We can express this responsiveness directly as a function of the induced posterior. Transforming the contracting problem in this fashion into belief space, it reads

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

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

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

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

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.*

From the perspective of the first period, the profit in the terminal period induces a continuation value

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

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 ( $\Delta a$ ). 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 ( $\Delta b$ ). It dominates and learning reduces profits whenever complementarities are sufficiently strong, i.e. when the technology is log-supermodular in effort and ability.<sup>18</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*

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

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 required 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 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

---

<sup>18</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, see Remark 1. A sufficient condition is that the probability of low output is log-supermodular in effort and ability.

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>19</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  $S = \{G, B\}$  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 may not hold when the wage function is fixed exogenously.

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

<sup>19</sup>To see this, note that  $\frac{w'''(x)}{w''(x)} = \frac{2\gamma-1}{1-\gamma} \frac{1}{x}$  for CRRA utility.

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,

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

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 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.

### 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 (2020). In the main text, we sketch the main steps of the argument, while we relegate the explicit duality arguments that are required to justify our approach to the appendix.

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,

$$\begin{aligned} \mathcal{L}(m, w; (\lambda_P, \lambda_{IC})) = & \int \left\{ P_\mu^1 Y + \delta \Pi_2(\hat{\mu}) - w(\hat{\mu}) \right. \\ & + \lambda_P (u(w(\hat{\mu})) - c - U) \\ & \left. + \lambda_{IC} \left( \frac{b + \Delta b \mu}{(\Delta a + \Delta b) \mu (1 - \mu)} (\hat{\mu} - \mu) u(w(\hat{\mu})) - c \right) \right\} d\hat{\mu}. \end{aligned} \quad (13)$$

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

**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 wage schedule follows from pointwise optimization of the Lagrangian.

$$w^*(\lambda, \hat{\mu}) = \max\left\{0, u'^{-1}\left(\lambda_P + \frac{b + \Delta b \mu}{(\Delta a + \Delta b) \mu (1 - \mu)} (\hat{\mu} - \mu)\right)\right\} \quad (14)$$

Plugging this function back into the Lagrangian of the problem, it is written as an expectation of a function of the posterior

$$\sup_w \mathcal{L}(m, w; \lambda) = \int \ell^*(\hat{\mu}; \lambda) m(\hat{\mu}) d\hat{\mu} \quad (15)$$

where  $\ell^*$  is the integrand of (13) evaluated at (14).

**Information Design** Therefore, the information design problem for a given  $\lambda$  is of standard form. The principal simply maximizes the expectation of a function of posteriors,

$$\sup_{m \text{ s.t. (BP)}} \int \ell^*(\hat{\mu}; \lambda) m(\hat{\mu}) d\hat{\mu}. \quad (16)$$

Such an information design problem can be solved via concavification. 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>20</sup> that

$$\begin{aligned} \frac{\partial^2}{\partial \hat{\mu}^2} \ell^*(\hat{\mu}; \lambda) = & \lambda_{IC}^2 \left[ \frac{b + \Delta b \mu}{(\Delta a + \Delta b) \mu (1 - \mu)} \right]^2 \rho' \left( \lambda_P + \lambda_{IC} \frac{b + \Delta b \mu}{(\Delta a + \Delta b) \mu (1 - \mu)} (\hat{\mu} - \mu) \right) \\ & + \delta \Pi_2''(\hat{\mu}) \end{aligned} \quad (17)$$

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 of (17) 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

---

<sup>20</sup>In the main text, we suppress boundary conditions related to the non-negativity constraint on wages.

continuation value. It is negative: the principal wants to keep the agent uninformed in order to reduce agency costs in the next period.

Furthermore, we have that

$$\begin{aligned} \frac{\partial^3}{\partial \hat{\mu}^3} \ell^*(\hat{\mu}; \lambda) = & \lambda_{IC}^3 \left[ \frac{b + \Delta b \mu}{(\Delta a + \Delta b) \mu (1 - \mu)} \right]^3 \rho''(\lambda_P + \lambda_{IC} \frac{b + \Delta b \mu}{(\Delta a + \Delta b) \mu (1 - \mu)} (\hat{\mu} - \mu)) \\ & + \delta \Pi_2'''(\hat{\mu}) > 0 \end{aligned} \quad (18)$$

The third derivative 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,  $\rho'' \geq 0$ .<sup>21</sup> The second term is determined by the shape of the continuation value. The principal is less information averse for high posteriors.

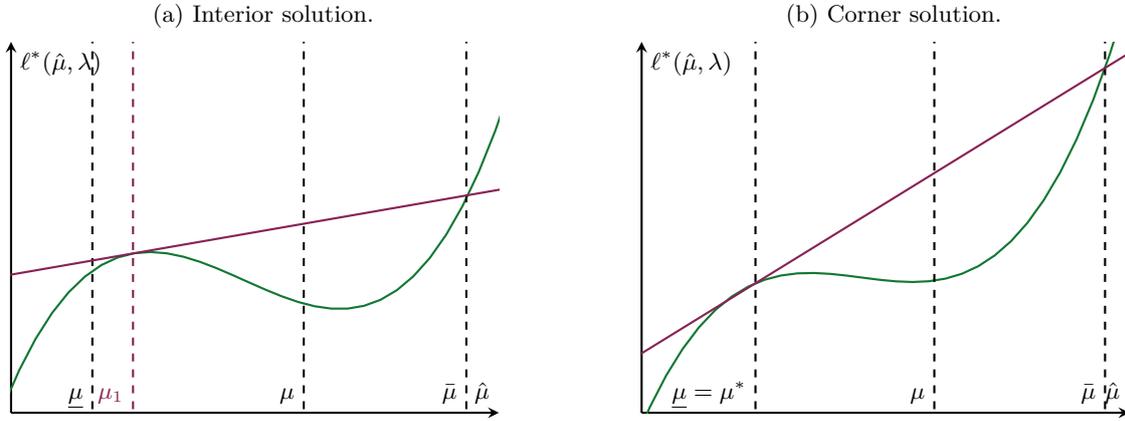


Figure 1: The concavification of  $\ell^*$  at  $\mu$ .

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  $\ell^*$  is concave for low posteriors and convex for high posteriors. The optimal information structure is fully informative at the top and uses partial pooling at the bottom (Fig. 1a). Finally, as  $\lambda_{IC}$  increases further,  $\ell^*$  becomes globally convex as the costs of incentives overwhelm the gains from concealing information. The resulting evaluation structure is fully informative (Fig. 1b).

The arguments in the appendix establish through a series of lemmas that a solution to the problem exists and is characterized by the two-step procedure above. QED.

<sup>21</sup>To see this, note that  $w''(u(x)) = -\frac{w''(x)}{[w'(x)]^3}$  and that  $\rho'(x) = -\frac{[u'(f(x))]^3}{u''(f(x))}$  for the strictly increasing  $f(x) = u'^{-1}(\frac{1}{x})$ . Hence, the former is decreasing (Assumption 1.3) if and only if the latter is increasing ( $\rho'' \geq 0$  as required).

This condition sufficient but I conjecture that the restriction to utility functions with  $\rho'' \geq 0$  is far from necessary. Instead, it is a result of the proof approach that requires establishing properties of the Lagrangian that are uniform across multipliers;  $\frac{\partial^3}{\partial \hat{\mu}^3} \ell^*(\hat{\mu}; \lambda)$  is positive for all multipliers only if  $\rho'' \geq 0$ .

### 3.4 Analysis of the Solution and Comparative Statics

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

$$\ell^*(\mu^*, \lambda(\mu^*)) + \frac{\partial \ell^*}{\partial \hat{\mu}} \Big|_{(\hat{\mu}, \lambda)=(\mu^*, \lambda(\mu^*))} (\bar{\mu} - \mu^*) = \ell(\bar{\mu}, \lambda(\mu^*)), \quad (19)$$

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 incentive compatibility constraint. For the graphical representation of our analysis this implies that, as we vary the tangent point in Figure 1 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 (19) into a more concrete form:

$$\frac{c^2}{2} \left( \frac{(\Delta a + \Delta b)\mu(1 - \mu)}{b + \Delta b\mu} \frac{\bar{\mu} - \mu^*}{(\bar{\mu} - \mu)(\mu - \mu^*)} \right)^2 = \delta (\Pi_2(\mu^*) + \Pi_2'(\bar{\mu})(\bar{\mu} - \mu^*) - \Pi_2'(\bar{\mu})) \quad (20)$$

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

$$\sigma(\mu^*) = \frac{1 - P_\mu \mu^* - \mu}{P_\mu \bar{\mu} - \mu^*} \in [0, 1], \quad (21)$$

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 4.** *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, and for  $u(x) = \sqrt{2x}$ , 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

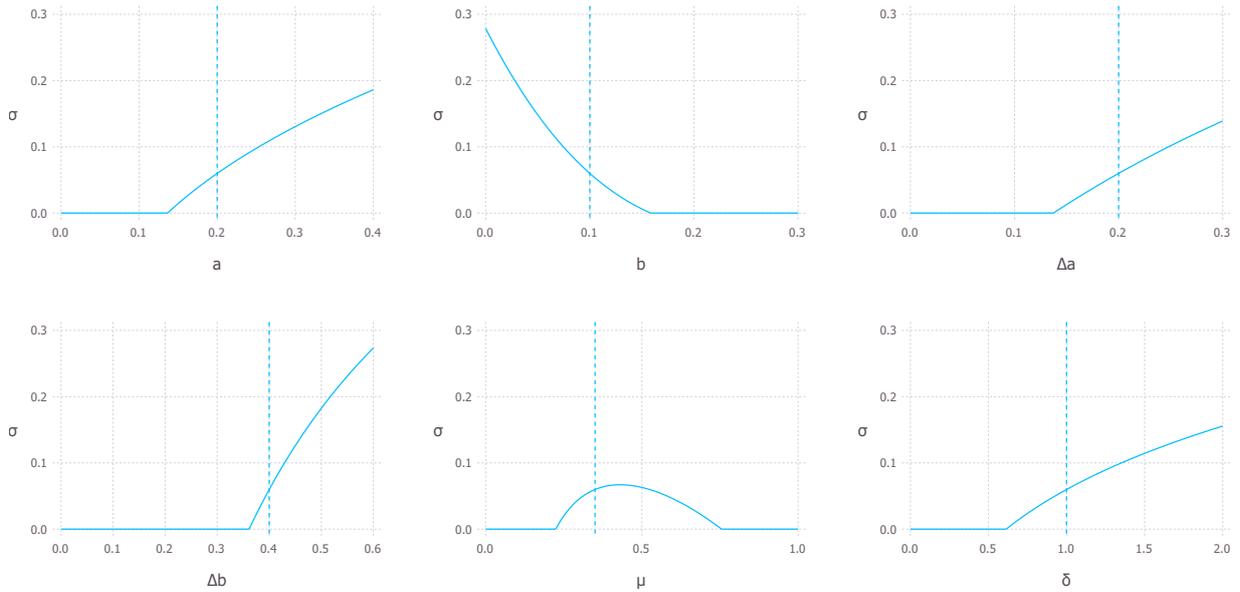


Figure 2: Comparative statics with  $u(x) = \sqrt{2x}$  around  $(a, b, \Delta a, \Delta b) = (0.2, 0.1, 0.2, 0.4)$ ,  $\mu = 0.35$ ,  $\delta = 1$ .

agency costs in that period more important, thus the evaluation structure becomes more informative. Numerical computations (Figure 2) show that the noise in the evaluation is highest when there is a lot of uncertainty about the agent's type (intermediate priors), when the impact of ability on output is high (both directly,  $\Delta a$ , and through the impact of effort,  $\Delta b$ ), and when high ability is relatively more essential for the impact of effort (low  $\Delta b$ ) and less essential for the baseline level of output (high  $a$ ).

*Remark 1* (Substitutes). If effort and ability enter the likelihood of high output linearly (i.e.  $\Delta b = 0$ ), agent learning is not costly and  $\Pi_2$  is convex. Therefore, the optimal evaluation is fully informative. If the two are sufficiently strong substitutes (i.e.  $-(1-a)\Delta b > b\Delta a$ ),  $\Pi_2$  is concave with  $\Pi_2''' < 0$  and one can show under a modification of Assumption 1 (replacing  $w''' \leq 0$  by  $w''' \geq 0$ ) that the optimal evaluation is binary and lenient. In this setting, workers who are too sure of themselves are disproportionately expensive to motivate and therefore it is optimal to reduce the informativeness of the good evaluation.

## 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>22</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>23</sup> Some individuals are also systematically underconfident and this trait is common in some demographic groups.<sup>24</sup>

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>25</sup>

#### 4.1 Contracting with Heterogeneous Priors

We solve the contracting problem with heterogeneous priors (de la Rosa, 2011; Dumav et al., 2021). 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>26</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.

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

---

<sup>22</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>23</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>24</sup>There is some evidence that women tend to be underconfident, for example (Niederle and Vesterlund, 2007; Hügelschäfer and Achtziger, 2014).

<sup>25</sup>Indeed, supporting students' self-esteem is often cited as a reason for grade inflation in schools and universities (Boretz, 2004).

<sup>26</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.

( $a\Delta b - b\Delta a > 0$ ) 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  $u(x) = \sqrt{2x}$ .

### 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 and the agent now have heterogeneous priors. In particular, the belief of the principal is degenerate and we therefore write the problem in terms of the posterior 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,

$$\int \hat{\mu} m(\hat{\mu}) d\hat{\mu} = \mu \quad (22)$$

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>27</sup> Then, the probability of  $\hat{\mu}(s)$  according to the agent is

$$m(\hat{\mu}(s)) = p(s|y_L) + (a + b + (\Delta a + \Delta b)\mu) [p(s|y_H) - p(s|y_L)] \quad (23)$$

According to the principal, this event has probability

$$\begin{aligned} m_P(\hat{\mu}(s)) &= p(s|y_L) + (a + b + (\Delta a + \Delta b)\eta) [p(s|y_H) - p(s|y_L)] \\ &= m(\hat{\mu}(s)) + (\eta - \mu)(\Delta a + \Delta b) [p(s|y_H) - p(s|y_L)] \\ &= \left[ \eta \frac{\hat{\mu}(s)}{\mu} + (1 - \eta) \frac{1 - \hat{\mu}(s)}{1 - \mu} \right] m(\hat{\mu}(s)) \end{aligned} \quad (24)$$

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>28</sup>

---

<sup>27</sup>In the optimal evaluation structure, there is at most one such signal by a straightforward extension of Proposition 1.

<sup>28</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

$$\Pi_t^\eta(\mu) = \max_{m,w} P_\eta^1 Y + \int (\delta \Pi_{t+1}^\eta(\hat{\mu}) - w(\hat{\mu})) D^\eta(\mu, \hat{\mu}) m(\hat{\mu}) d\hat{\mu} \quad (25)$$

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

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

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

## 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.

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

$$D^\eta(\mu, \hat{\mu}) \Pi_2^\eta(\hat{\mu}) \quad (26)$$

Learning now has two effects. First, it creates a dispersion of the agent's posterior. This is costly for the principal, since  $\frac{\partial^2}{\partial \bar{\mu}^2} \Pi_2^\eta(\hat{\mu}) < 0$  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, 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. Learning reduces profits with optimism and has an ambiguous impact with pessimism. Let us summarize the preceding discussion.

**Proposition 5.** *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.*

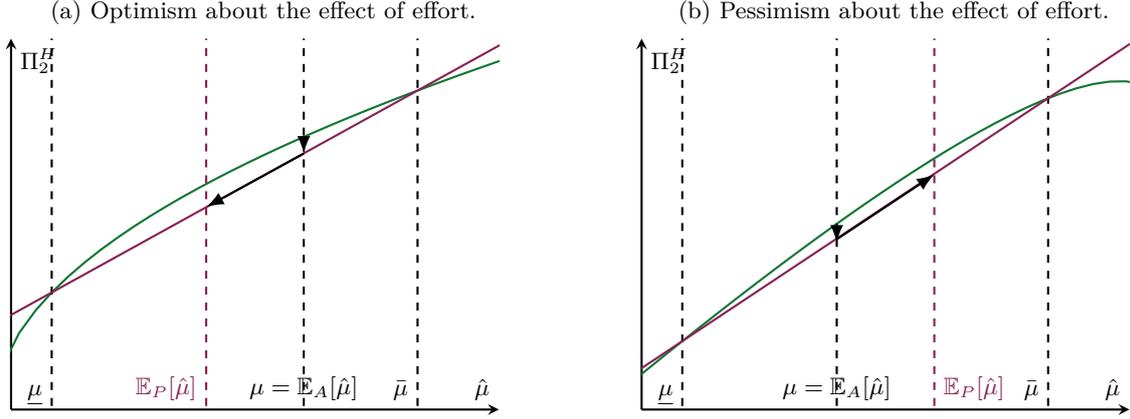


Figure 3: The effect of information on the principal's continuation value.

The impact of information is determined by

$$\frac{\partial^2}{\partial \hat{\mu}^2} (D^\eta(\mu, \hat{\mu}) \Pi_2^\eta(\hat{\mu})) = \underbrace{\left[ \eta \frac{\hat{\mu}}{\mu} + (1 - \eta) \frac{1 - \hat{\mu}}{1 - \mu} \right] \Pi_2^{\eta''}(\hat{\mu})}_{\text{dispersion}} + 2 \underbrace{\frac{\eta - \mu}{(1 - \mu)\mu} \Pi_2^{\eta'}(\hat{\mu})}_{\text{drift}} \quad (27)$$

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

### 4.3 Initial Period

The shape of the optimal evaluation structure is determined by two factors: First, based on the 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

$$\frac{\partial^3}{\partial \hat{\mu}^3} (D^\eta(\mu, \hat{\mu}) \Pi_2^\eta(\hat{\mu})) = \underbrace{\left[ \eta \frac{\hat{\mu}}{\mu} + (1 - \eta) \frac{1 - \hat{\mu}}{1 - \mu} \right] \Pi_2^{\eta'''}(\hat{\mu})}_{\text{dispersion}} + 3 \underbrace{\frac{\eta - \mu}{(1 - \mu)\mu} \Pi_2^{\eta''}(\hat{\mu})}_{\text{drift}} \quad (28)$$

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.*

$$p(G|y_H) = 1 - \sigma, \quad p(B|y_H) = \sigma, \quad p(G|y_L) = 0 \quad p(B|y_L) = 1. \quad (29)$$

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

$$p(G|y_H) = 1, \quad p(B|y_H) = 0, \quad p(G|y_L) = \sigma \quad p(B|y_L) = 1 - \sigma. \quad (30)$$

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 has an incentive to do so.

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>29</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,  $(\mu_P, \hat{\mu}) \rightarrow (w_L, w_H)(\mu_P, \hat{\mu})$ , and (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. We say that a PBE satisfies *no-holdup* if the agent's participation constraint is binding in almost all on-path second period contract offers.<sup>30</sup> 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>31</sup>

The outcome of Theorem 1 is achieved as the unique equilibrium in this class.

*Remark 2.* 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 no-holdup 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 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>32</sup> To see that this equilibrium is

<sup>29</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 we study can be extended naturally.

<sup>30</sup>Without such a refinement, the equilibrium could grant intertemporal commitment. This would allow the principal to smooth out bonus payments across periods, yielding higher profits through a channel orthogonal to the acquisition of private information.

<sup>31</sup>Orlov et al. (2020) 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).

<sup>32</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

principal preferred, note that in any no-holdup 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 by using it as the basis of incentives in the first period, thus reducing agency costs.

The principal prefers to commit not to reveal additional information through the contract offer. Passive beliefs effectively 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.

*Remark 3.* When the information structure is observed, any equilibrium<sup>33</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 share a common belief 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 deviation to  $e_1 = 0$ , however, the agent updates his beliefs according to

$$\tilde{\mu}(s) = \mu \frac{p(s|y_L) + (a + \Delta a) [p(s|y_H) - p(s|y_L)]}{p(s|y_L) + (a + \mu\Delta a) [p(s|y_H) - p(s|y_L)]} \quad (31)$$

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>34</sup> A deviation in the first period is more profitable for the agent because of this belief-manipulation effect.<sup>35</sup> We now analyze this model,

---

not the case in our game. There are also no proper subgames to which they could be applied. [Ekmekci and Kos \(2021\)](#) 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.

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>33</sup>Among no-holdup 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  $m_P(\cdot, \hat{\mu})$ .

<sup>34</sup>This assumes that the principal does not elicit the agent's belief at the beginning of the second period. In such a mechanism, however, truth-telling would need to be preferable to imitating the type that realizes on path. Hence, a screening mechanism in the second period cannot reduce the post-deviation payoff and therefore does not affect the optimal contract.

<sup>35</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

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. 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 first period IC reads

$$\int_S (p(s|y_L) + (a + b + \mu\Delta b) [p(s|y_H) - p(s|y_L)]) [w(s) + U] ds - c \geq \int_S (p(s|y_L) + a [p(s|y_H) - p(s|y_L)]) \cdot \left[ w(s) + \max \left\{ w_L(\hat{\mu}(s)) + P_\mu^1 \frac{c}{b + \Delta b \hat{\mu}(s)} - c, w_L(\hat{\mu}(s)) + P_\mu^0 \frac{c}{b + \Delta b \hat{\mu}(s)} \right\} \right] ds \quad (32)$$

The 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 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 in 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 it as

$$\int \left\{ \frac{(b + \mu\Delta b)}{\mu(1 - \mu)\Delta b} (\hat{\mu} - \mu) u(w(\hat{\mu})) - \left[ 1 - \frac{(b + \mu\Delta b)}{\mu(1 - \mu)\Delta b} (\hat{\mu} - \mu) \right] \max\{0, c\Delta b \frac{\mu - \hat{\mu}}{b + \hat{\mu}\Delta b}\} \right\} m(\hat{\mu}) d\hat{\mu} \geq c \quad (33)$$

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.

*Remark 4.* 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

---

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.

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 in Theorem 1. The kink in the IC constraint, however, raises the possibility of a third, uninformative signal. This signal can help to economize on the costs caused by belief-manipulation (33). 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 costlessly 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 only at the end of the employment relationship and that doing so is costless e.g. because the agent only consumes at the end of the second period. Then, informative wages do not lead to learning and hence using a fully informative evaluation is optimal. Any feature of the model, however, that makes it costly or impossible to delay informative incentive payments reinstates the trade-off between learning and incentives analyzed in this paper. 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 problem is equivalent to period-by-period contracting. Noisy and (weakly) tough performance evaluation is again optimal and this extends by continuity to interior discount rates. Also simple risk aversion implies that it is costly to delay informative incentives, as it would be optimal to smooth out bonus payments across both periods.

Our results continue to hold if the principal can postpone only payments, but not information. Suppose that the first period contract specifies not only a wage this period, but also a continuation value. Our results generalize to this model.

*Remark 5.* Suppose that  $u(w) = \sqrt{2w}$  and that in the first period, the principal can commit to signal-contingent wages and continuation values. The optimal information structure is essentially unique, binary, and (weakly) tough.

The case with full commitment raises considerable difficulties and is beyond the scope of this paper. This is because of the interaction between full commitment and the belief-manipulation problem. The dynamic contract cannot condition on the true effort exerted in the first period, as this would resolve the moral hazard problem. Therefore, the full commitment contract has to deal with the dynamic constraint (33) outlined in previous section. As a result, the principal may find it optimal to commit to excessive bonus

payments in the final period to relax this constraint by inducing a learning motive in the agent. To analyze this problem, the contracts in both periods need to be designed jointly with the information structure, which is intractable.<sup>36</sup>

## 5.4 Long-Run Simulations

In the main sections, we analyze the twice-repeated problem. Consider now the same period-by-period problem, but repeated for  $T > 2$  periods. Characterizing the solution of this problem analytically is not tractable. Numerical analysis, however, shows that the results from the two-period problem generalize: The optimal evaluation is binary and tough in all periods. The longer the remaining time horizon, the more important the dynamic learning channel. Therefore, the optimal evaluation is less precise in early periods and becomes more precise over time (Figure 4). The worker is left in the dark through a noisy and tough evaluation early in the relationship when additional information affects many future incentive compatibility constraints. There is significant uncertainty under the optimal evaluation even after ten periods, with fewer very low and more moderately high posterior beliefs compared to a fully informative one (Figure 5b)

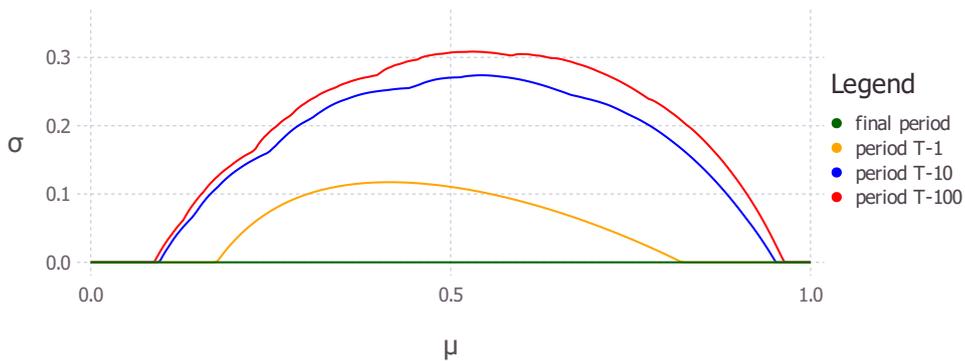


Figure 4: The probability of a false negative evaluation decreases over time for any prior.

## 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.

<sup>36</sup>One possible work-around is to consider the problem with full commitment when the contract terms in the second period can condition on true effort in the first period, the expected utility of the agent, however, is independent of this information conditional on the evaluations. This problem is akin to commitment to a continuation value and our results continue to hold as in Remark 5.

(a) The evolution of beliefs under the dynamically optimal evaluation. (b) The CDF at  $T = 10$  (posteriors on the  $y$ -axis).

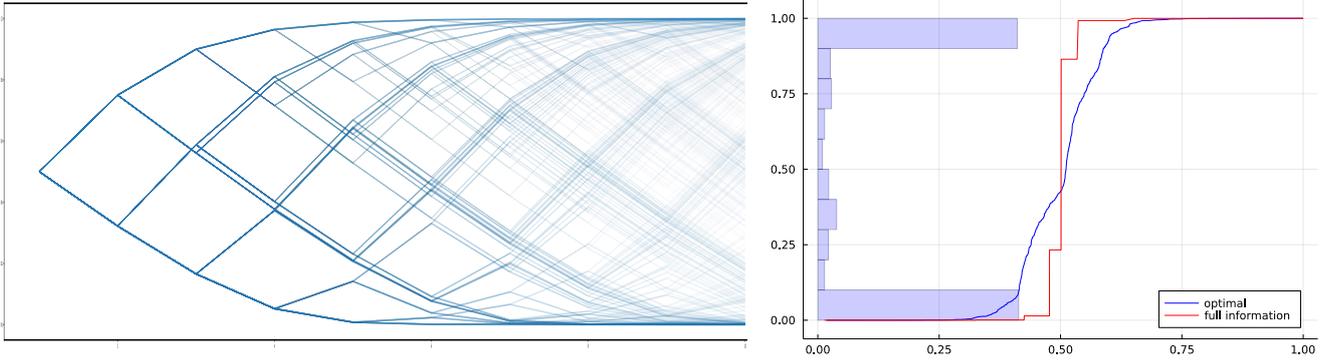


Figure 5: Belief dynamics in the Problem with  $T = 10$  (parameters as in Fig. 2)

When effort and ability are log-complements, 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 in that sense part of the optimal organization design.

Our results inform not only the optimal evaluation of employee performance, but are also suggestive of the selection of information sources. Among measures that combine information about effort and ability, the principal prefers measure that are less sensitive to ability. This may motivate secrecy about salary differences between employees in the same job as opposed to differences in bonuses, as base-pay reflects the principal's estimate of ability while bonuses more directly reward effort. Monitoring effort itself and the conduct of employees more generally instead of output remains desirable. These patterns are in sharp contrast with models of implicit incentives through career concerns. In such models, the fact that a signal combines information about effort and ability is not a friction but 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 is an interesting avenue for future research.

## A Proofs

Note that we use the letter  $w$  both to denote the promised wage as a function of signals  $s$  and posteriors  $\hat{\mu}$ , as well as for  $w = u^{-1}$  mapping from promised utilities to the monetary cost of providing it. Which function is denoted in each specific instance is apparent from its arguments and context.

*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

$$\begin{aligned} \max_{m,v} & P_\mu Y + \int_S \left( \delta \Pi_{t+1}(\hat{\mu}(s)) - w(v(s)) \right) m(s) ds \\ \text{s.t.} & \int_S v(s) m(s) ds - c \geq U \\ & \int_S \left( b + \Delta b \mu \right) \frac{\mu(s) - \mu}{\mu(1 - \mu)(\Delta a + \Delta b)} v(s) m(s) ds \geq c \quad (\text{IC}) \\ & \int \hat{\mu} m(\hat{\mu}) d\hat{\mu} = \mu; \quad \text{supp}(m) \subset [\underline{\mu}, \bar{\mu}] \quad (\text{BP}) \end{aligned}$$

where  $v(s)$  is the promised utility at signal  $s$  and where we used the representation of the IC in (6) while writing the posterior as a function of  $s$ . 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  $p(s|y_L) = 0$ , which attains the upper bound, and minimized for  $p(s|y_H) = 0$ , which attains the lower bound.  $\square$

*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  $\frac{\partial}{\partial \mu} P_\mu \frac{c}{b + \Delta b \mu} \propto \Delta a b - \Delta b a < 0$ , this is implied by  $U > \frac{a+b}{b} c$ . It is easy to

verify that

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

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 then follows.  $\square$

*Proof of Theorem 1:* The contracting problem is equivalent to

$$\sup_{w, m \text{ s.t. (BP)}} \inf_{\lambda \geq 0} \mathcal{L}(m, w; (\lambda_P, \lambda_{IC})). \quad (34)$$

To see this, note that as

$$\inf_{\lambda \geq 0} \mathcal{L}(m, w; \lambda) = \begin{cases} \int m(\hat{\mu}) (P_\mu^1 Y + \delta \Pi_2(\hat{\mu}) - w(\hat{\mu})) d\hat{\mu} & \text{if (P)\&(IC) are satisfied} \\ -\infty & \text{else} \end{cases}, \quad (35)$$

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*.

Fix a distribution  $m$  satisfying (BP) and consider the problem of finding optimal wages subject to the participation and incentive constraint. It can be solved by point-wise optimization, arriving at (14).

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

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

*Proof of Lemma:* If  $m$  is degenerate (a point mass on  $\mu$ ), the problem is infeasible and hence both the primal and dual value are  $-\infty$ . If  $m$  is nondegenerate, it is easy to see that the problem can be written in utility space where the objective is a concave functional and the constraints are linear. Furthermore, a strictly feasible utility promise exists (e.g., pay  $U$  after every posterior with a suitable large bonus if and only if  $\hat{\mu} > \mu$ ). Therefore, by standard results (e.g. [Luenberger, 1969](#), p. 224), the problem satisfies strong duality. Therefore, so does the wage setting problem since the two are equivalent.  $\triangle$

Consider now the Lagrangian that results from plugging in for  $w$  from (14).

$$\sup_w \mathcal{L}(m, w; \lambda) = \int \ell^*(\hat{\mu}; \lambda) m(\hat{\mu}) d\hat{\mu} \quad (36)$$

The information design problem given  $\lambda$  reads

$$\sup_{m \text{ s.t. (BP)}} \int \ell^*(\hat{\mu}; \lambda) m(\hat{\mu}) d\hat{\mu}, \quad (37)$$

and can therefore be solved via concavification of  $\ell^*$ .

**Lemma 2.** *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  $m(\mu^*) \in [\frac{\mu - \underline{\mu}}{\bar{\mu} - \underline{\mu}}, 1]$ .*

*Proof of Lemma:* From (17), it is easy to see that

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

at an interior wage and  $\frac{\partial^3}{\partial \bar{\mu}^3} \ell^*(\hat{\mu}; \lambda) = \delta \Pi_2'''(\hat{\mu})$  when the wage is zero.<sup>37</sup> It is elementary but tedious to show that

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

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  $\frac{\partial^3}{\partial \bar{\mu}^3} \ell^*(\hat{\mu}; \lambda) > 0$  for all  $\hat{\mu}$  and  $\lambda$ .

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  $f$  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  $\mu$ ,

$$\Psi(\lambda_P, \lambda_{IC}) := \{\psi \in [\underline{\mu}, \bar{\mu}] | \exists \psi' \in [\underline{\mu}, \bar{\mu}], \alpha \in [0, 1] \text{ s.t. } \alpha \psi + (1 - \alpha) \psi' = \mu \text{ and} \quad (38)$$

$$\text{cav} \ell^*(\mu; \lambda_P, \lambda_{IC}) = \alpha \ell^*(\psi; \lambda_P, \lambda_{IC}) + (1 - \alpha) \ell^*(\psi'; \lambda_P, \lambda_{IC}) \quad (39)$$

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  $\mu$  (since  $\ell^{*'''} > 0$ ) 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  $M(\mu^*) + M(\bar{\mu}) = 1$ .  $\triangle$

Applying this result to the original problem requires another step of duality.

**Lemma 3.** *The information design problem satisfies strong duality, i.e.*

$$\sup_{m \text{ s.t. (BP)}} \inf_{\lambda \geq 0} \int \ell^*(\hat{\mu}; \lambda) m(\hat{\mu}) d\hat{\mu} = \inf_{\lambda \geq 0} \sup_{m \text{ s.t. (BP)}} \int \ell^*(\hat{\mu}; \lambda) m(\hat{\mu}) d\hat{\mu}. \quad (40)$$

*Proof of Lemma:* 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

<sup>37</sup>Note that  $\ell^{*'}$  is continuous at the nonnegativity constraint as the wage is locally zero and that  $\ell^{*''}$  jumps upwards as  $\rho' \geq 0$ .

linear in  $m$ . Continuity in  $\lambda$  is immediate. To see quasi-convexity in  $\lambda$ , note that by an envelope argument

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

and similarly for  $\lambda_{IC}$  and hence the Hessian of  $\int \ell^*(\hat{\mu}; \lambda) m(\hat{\mu}) d\hat{\mu}$  is given by

$$\begin{pmatrix} \int f(\hat{\mu}) d\hat{\mu} & \int f(\hat{\mu}) \frac{b + \Delta b \mu}{(\Delta a + \Delta b) \mu (1 - \mu)} (\hat{\mu} - \mu) d\hat{\mu} \\ \int f(\hat{\mu}) \frac{b + \Delta b \mu}{(\Delta a + \Delta b) \mu (1 - \mu)} (\hat{\mu} - \mu) d\hat{\mu} & \int f(\hat{\mu}) \left[ \frac{b + \Delta b \mu}{(\Delta a + \Delta b) \mu (1 - \mu)} (\hat{\mu} - \mu) \right]^2 d\hat{\mu} \end{pmatrix} \quad (41)$$

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  $\int f(\hat{\mu}) g_1(\hat{\mu}) g_2(\hat{\mu}) d\hat{\mu}$  defines an inner product  $\langle \cdot, \cdot \rangle_f$  (between functions that share support with  $m$ ), and  $\int \ell^*(\hat{\mu}; \lambda) m(\hat{\mu}) d\hat{\mu}$  is weakly convex by Cauchy-Schwarz, as the determinant of the Hessian reads

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

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

Therefore, the problem satisfies the conditions of Sion's Minimax Theorem and we have

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

Using Lemma 1 and 3, we have

$$\Pi_1(\mu) = \sup_{w, m \text{ s.t. (BP)}} \inf_{\lambda \geq 0} \mathcal{L}(m, w; \lambda) = \sup_{m \text{ s.t. (BP)}} \inf_{\lambda \geq 0} \sup_w \mathcal{L}(m, w; \lambda) = \inf_{\lambda \geq 0} \sup_{w, m \text{ s.t. (BP)}} \mathcal{L}(m, w; \lambda).$$

We can therefore simplify the general problem (7) using the properties of optimal evaluation structures from Lemma 2, i.e. we can restrict attention to binary information structures where the good signal only realizes after high output. This simplified problem is

$$\begin{aligned} \max_{\mu^*, m^*, w_l, w_h} & P_\mu^1 Y + m^* [\delta \Pi_2(\mu^*) - w_l] + (1 - m^*) [\delta \Pi_2(\bar{\mu}) - w_h] & (42) \\ \text{s.t.} & m^* u(w_l) + (1 - m^*) u(w_h) - c \geq U & (\text{P}^S) \\ & \frac{b + \Delta b \mu}{(\Delta a + \Delta b) \mu (1 - \mu)} [m^* (\mu^* - \mu) u(w_l) + (1 - m^*) (\bar{\mu} - \mu) u(w_h)] \geq c & (\text{IC}^S) \\ & m^* \mu^* + (1 - m^*) \bar{\mu} = \mu; \quad \mu^* \in [\underline{\mu}, \mu] & (\text{BP}^S) \end{aligned}$$

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 (42) has a unique solution. The optimal information structure is non-degenerate ( $\mu^* < \mu$ ).*

*Proof of Lemma:* To show nondegeneracy in the simplified problem, we need to show that the optimal distribution of posteriors is nondegenerate. To this purpose, we show that there exists an  $\bar{\epsilon}_1$  such that  $\mu^* \leq \mu - \bar{\epsilon}_1$ . 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

$$m(\mu^*)\mu^* + m(\bar{\mu})\bar{\mu} = \mu$$

$$m(\bar{\mu}) = \frac{\mu - m(\mu^*)\mu^*}{\bar{\mu} - \mu} = \frac{\mu - m(\mu^*)\mu - \epsilon}{\bar{\mu} - \mu} = \epsilon \frac{m(\mu^*)}{\bar{\mu} - \mu} \leq \epsilon \frac{1}{\bar{\mu} - \mu}$$

In the IC constraint, we have

$$c \leq \frac{b + \Delta b\mu}{(\Delta a + \Delta b)\mu(1 - \mu)} [m^*(\mu^* - \mu)u(w^*) + (1 - m^*)(\bar{\mu} - \mu)u(\bar{w})]$$

$$\leq \frac{b + \Delta b\mu}{(\Delta a + \Delta b)\mu(1 - \mu)} \left[ \epsilon \frac{1}{\bar{\mu} - \mu} (\bar{\mu} - \mu)(u(\bar{w}) - u(\underline{w})) \right]$$

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  $\leq c_1 - \epsilon \cdot c_2 w(\epsilon^{-1}) \rightarrow -\infty$  for suitable constants, which is clearly not optimal.

Hence, the optimal distribution of posteriors is nondegenerate. It is easy to see that the constraints have to be binding and consequently the wages can be expressed from the constraints. Therefore, a solution exists. Is is unique since the problem is concave with a convex constraint sets.  $\triangle$

Since a wage function and a distribution over posteriors solve the original problem if and only if they induce a solution in the simplified problem, this concludes the proof. The contract is unique in the sense that the signal structure is unique up to duplication and the wage function is determined on all signals that realize with a positive probability.  $\square$

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

$$\frac{b + \Delta b\mu}{(\Delta a + \Delta b)\mu(1 - \mu)} [m^*(\mu^* - \mu)u(w^*) + (1 - m^*)(\bar{\mu} - \mu)u(\bar{w})] =$$

$$\frac{b + \Delta b\mu}{(\Delta a + \Delta b)\mu(1 - \mu)} \left[ \frac{\bar{\mu} - \mu}{\bar{\mu} - \mu^*} (\mu^* - \mu)u(w^*) + \frac{\mu - \mu^*}{\bar{\mu} - \mu^*} (\bar{\mu} - \mu)u(\bar{w}) \right] =$$

$$\frac{b + \Delta b\mu}{\Delta b\mu(1 - \mu)} \frac{(\bar{\mu} - \mu)(\mu - \mu^*)}{\bar{\mu} - \mu^*} [u(\bar{w}) - u(w^*)] \geq c$$

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)

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

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

$$\delta [\Pi_2(\mu^*) - \Pi_2(\bar{\mu}) + (\bar{\mu} - \mu^*) \Pi_2'(\mu^*)] = \frac{1}{2} \left( \frac{(\Delta a + \Delta b) \mu (1 - \mu)}{b + \Delta b \mu} c \right)^2 \frac{(\bar{\mu} - \mu^*)^2}{(\bar{\mu} - \mu)^2 (\mu - \mu^*)^2} \quad (43)$$

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

*Proof of Proposition 4:* The comparative statics follow from implicitly differentiating (43). First, note that LHS is decreasing in  $\mu^*$  (this follows immediately from the concavity of  $\Pi_2$ ) and that the RHS is increasing in  $\mu^*$  (as  $\bar{\mu} > \mu$ ). Therefore

$$\frac{d\mu^*}{d\delta} = \frac{[\Pi_2(\mu^*) - \Pi_2(\bar{\mu}) + (\bar{\mu} - \mu^*) \Pi_2'(\mu^*)]}{\frac{d}{d\mu^*} \left[ \frac{1}{2} \left( \frac{(\Delta a + \Delta b) \mu (1 - \mu)}{b + \Delta b \mu} c \right)^2 \frac{(\bar{\mu} - \mu^*)^2}{(\bar{\mu} - \mu)^2 (\mu - \mu^*)^2} \right] - \frac{d}{d\mu^*} \delta [\Pi_2(\mu^*) - \Pi_2(\bar{\mu}) + (\bar{\mu} - \mu^*) \Pi_2'(\mu^*)]} > 0$$

and hence  $\frac{d\sigma}{d\delta} > 0$ . To see the statement about costs, let  $c_t$  denote the cost of effort in period  $t$ . Then

$$\frac{d\mu^*}{dc_1} = \frac{c_1 \left( \frac{(\Delta a + \Delta b) \mu (1 - \mu)}{b + \Delta b \mu} \right)^2 \frac{(\bar{\mu} - \mu^*)^2}{(\bar{\mu} - \mu)^2 (\mu - \mu^*)^2}}{\frac{d}{d\mu^*} \delta [\Pi_2(\mu^*) - \Pi_2(\bar{\mu}) + (\bar{\mu} - \mu^*) \Pi_2'(\mu^*)] - \frac{d}{d\mu^*} \left[ \frac{1}{2} \left( \frac{(\Delta a + \Delta b) \mu (1 - \mu)}{b + \Delta b \mu} c \right)^2 \frac{(\bar{\mu} - \mu^*)^2}{(\bar{\mu} - \mu)^2 (\mu - \mu^*)^2} \right]} < 0$$

and hence  $\frac{d\sigma}{dc_1} < 0$ . To see that

$$\frac{d\mu^*}{dc_2} = \frac{\delta \frac{d}{dc_2} [\Pi_2(\mu^*) - \Pi_2(\bar{\mu}) + (\bar{\mu} - \mu^*) \Pi_2'(\mu^*)]}{\frac{d}{d\mu^*} \left[ \frac{1}{2} \left( \frac{(\Delta a + \Delta b) \mu (1 - \mu)}{b + \Delta b \mu} c \right)^2 \frac{(\bar{\mu} - \mu^*)^2}{(\bar{\mu} - \mu)^2 (\mu - \mu^*)^2} \right] - \frac{d}{d\mu^*} \delta [\Pi_2(\mu^*) - \Pi_2(\bar{\mu}) + (\bar{\mu} - \mu^*) \Pi_2'(\mu^*)]} > 0$$

and hence  $\frac{d\sigma}{dc_2} < 0$ , it remains to show that  $\frac{d}{dc_2} [\Pi_2(\mu^*) - \Pi_2(\bar{\mu}) + (\bar{\mu} - \mu^*) \Pi_2'(\mu^*)] > 0$ . This equation is only tractable for  $u(w) = \sqrt{2w}$ . Then, we have

$$\begin{aligned} & \frac{d}{dc_2} [\Pi_2(\mu^*) - \Pi_2(\bar{\mu}) + (\bar{\mu} - \mu^*) \Pi_2'(\mu^*)] = \\ & = \frac{d}{dc_2} c_2^2 \left[ \frac{(1 - P_{\bar{\mu}}) P_{\bar{\mu}}}{(b + \Delta b \bar{\mu})^2} - \frac{(1 - P_{\mu^*}) P_{\mu^*}}{(b + \Delta b \mu^*)^2} + (\bar{\mu} - \mu^*) \frac{\Delta b (b + \Delta b \mu^* + 2a(1 - P_{\mu^*}))}{(b + \Delta b \mu^*)^3} \right] \\ & = 2c_2 \left[ \frac{(1 - P_{\bar{\mu}}) P_{\bar{\mu}}}{(b + \Delta b \bar{\mu})^2} - \frac{(1 - P_{\mu^*}) P_{\mu^*}}{(b + \Delta b \mu^*)^2} + (\bar{\mu} - \mu^*) \frac{\Delta b (b + \Delta b \mu^* + 2a(1 - P_{\mu^*}))}{(b + \Delta b \mu^*)^3} \right] > 0 \end{aligned}$$

Finally, from this it is apparent that both sides of (43) are proportional to  $c_1^2$  and  $c_2^2$ , respectively, if  $u(w) = \sqrt{2w}$ , which implies that  $\sigma$  is invariant to common changes in costs in that case.  $\square$

*Proof of Proposition 5:* To see the first statement, consider the problem in signal/utility space. Then, the cost of incentives is

$$\int ((a + b + \eta(\Delta a + \Delta b))p(s|y_H) + (1 - a - b - \eta(\Delta a + \Delta b))p(s|y_L)) w(v(s)) ds$$

and the constraints depend on

$$\int ((a + b + \mu(\Delta a + \Delta b))p(s|y_H) + (1 - a - b - \mu(\Delta a + \Delta b))p(s|y_L)) v(s) ds$$

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

To see monotonicity and concavity in the case of underconfidence, 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 follow from analogous straightforward but tedious computation.

To see the result on information, note that information corresponds to a mean preserving spread 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  $\eta = 1$ : The sign is ambiguous,

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

In particular, the expression is increasing in  $U$ . To see that the sign is truly ambiguous, note that

- if  $b > 0$ , the principal is information loving if  $U$  is sufficiently large, 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  $b < \bar{b}$ .

To see this, let  $U = \frac{a+b}{b}c - \delta$  in order to ensure that the nonnegativity 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.  $\square$

*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 (25) 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 2 goes through. We can apply the proof of Lemmas 4 and 3 mutatis mutandis and arrive at the Theorem.

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

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

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.  $\square$

## A.1 Private Information of the Principal

Consider the game as described in the text.

First, consider any weak PBE satisfying no-holdup. We will show that the principal profit is smaller than  $\Pi^*$ . On path, it induces a distribution over agent posteriors  $m(\hat{\mu})$  and conditional on the posterior  $\hat{\mu}$  a distribution over information structures and wage schedules. As the participation and incentive compatibility constraints are satisfied conditional on the agent's second stage information set, they are satisfied conditional on  $\hat{\mu}$ . Hence, this distribution over information structures and wage schedules satisfies the constraints of the second period problem for  $\hat{\mu}$ . By the proof of Proposition 5 above, the optimal contract is binary and independent of the principal's belief. Therefore, the continuation profit  $\Pi_2^{EQ}$  satisfies  $\int \Pi_2^{EQ}(\mu_p, \hat{\mu}) dm(\mu_P | \hat{\mu}) \leq \int \Pi_2^*(\mu_p, \hat{\mu}) dm(\mu_P | \hat{\mu}) = \Pi_2^*(\hat{\mu})$ . The principal's continuation value is dominated by that under the optimal contract. Similarly, in the first period, the equilibrium induces a distribution over wages and agent posteriors on-path that satisfies the conditions of the first period problem (7).<sup>38</sup> This implies that the first-period profit under the equilibrium is dominated by that under the optimal contract, as we set out to argue.

Second, consider a weak PBE with passive beliefs. Formally, we require that the agent does not update his beliefs about  $\theta$  based on the contract offer in either period. We will show that any such PBE induces a joint distribution over agent beliefs and wages that is identical to the one induced by the optimal contract up to a set of measure zero. Consider the contract offer stage in the second period. If the agent has posterior belief  $\hat{\mu}$ , the principal can achieve the optimal second period profit if and only if (up to inessential modifications of the contract) she offers the optimal contract characterized in the proof of Proposition 5 above. Therefore, the principal offers this contract in any such PBE. Consequently, by the martingale property, the principal's value of inducing posterior beliefs  $\mu_p, \hat{\mu}$  is  $\Pi_2^*(\hat{\mu})$ . In the first period, the agent has belief  $\mu$  on path (by Bayes' rule) and off-path (by passive beliefs). Therefore, the principal's best response is the solution to the contracting problem. Therefore, the equilibrium is outcome equivalent to the optimal contract.

Third, suppose that the agent observes the principal's information structure and the agent's posterior satisfies the restriction of FN 33. We will show that any such PBE induces a joint distribution over agent beliefs and wages that is identical to the one induced by the optimal contract up to a set of measure zero. Note that the principal can achieve  $\Pi^*$  by offering the optimal first-period contract and committing not to acquire private information. Furthermore, she can achieve no higher profit by the previous remark. By the uniqueness of the optimal contract (up to

<sup>38</sup>In a weak no-holdup PBE the agent may be misguided about the contract offered in the second period after a hypothetical deviation on the first period. Such beliefs can only strengthen the first-period IC constraint: On-path, the agent is held to the participation constraint (no-holdup); after a deviation, he might obtain a positive continuation surplus.

measure zero events), the equilibrium has to induce this outcome, otherwise the principal would obtain a strictly lower profit.

## A.2 Private Effort Choice

First, we provide details for the rewriting of the dynamic IC constraint (32). It is immediate that we can write

$$\begin{aligned} & \int_S (p(s|y_L) + (a + b + \mu\Delta b) [p(s|y_H) - p(s|y_L)]) [w(s) + U] ds - c \geq \\ & \int_S (p(s|y_L) + a [p(s|y_H) - p(s|y_L)]) [w(s) + U] ds + \\ & \int_S (p(s|y_L) + a [p(s|y_H) - p(s|y_L)]) \max \left\{ w_L(\hat{\mu}(s)) + P_\mu^1 \frac{c}{b + \Delta b \hat{\mu}(s)} - c - U, w_L(\hat{\mu}(s)) + P_\mu^0 \frac{c}{b + \Delta b \hat{\mu}(s)} - U \right\} ds \end{aligned} \quad (44)$$

Using the usual rewriting of the signal probabilities and noting that  $U = w_L(\hat{\mu}(s)) + P_\mu^1 \frac{c}{b + \Delta b \hat{\mu}(s)} - c$ , we can write

$$\begin{aligned} & \int m(\hat{\mu}) \frac{(b + \mu\Delta b)}{\mu(1 - \mu)\Delta b} (\hat{\mu} - \mu) w(\hat{\mu}) d\hat{\mu} - c \geq \\ & \int m(\hat{\mu}) \left( 1 - \frac{(b + \mu\Delta b)}{\mu(1 - \mu)\Delta b} (\hat{\mu} - \mu) \right) \max \left\{ (\mu - \hat{\mu})\Delta b \frac{c}{b + \Delta b \hat{\mu}}, 0 \right\} d\hat{\mu} \end{aligned} \quad (45)$$

which is the form given in (33). Note that the maximum is equal to zero iff  $\hat{\mu} \geq \mu$  and that wage setting is unaffected by this additional term. The partially maxed out Lagrangian reads

$$\int \tilde{\ell}(\hat{\mu}; \lambda) m(\hat{\mu}) d\hat{\mu}.$$

with

$$\tilde{\ell}(\hat{\mu}; \lambda) = \begin{cases} \ell^*(\hat{\mu}; \lambda) - \lambda_{IC} \left( 1 - \frac{(b + \mu\Delta b)}{\mu(1 - \mu)\Delta b} (\hat{\mu} - \mu) \right) (\mu - \hat{\mu})\Delta b \frac{c}{b + \Delta b \hat{\mu}} & \hat{\mu} \leq \mu \\ \ell^*(\hat{\mu}; \lambda) & \hat{\mu} > \mu \end{cases}$$

We have

$$\frac{\partial^3}{\partial \hat{\mu}^3} \left( -\lambda_{IC} \left( 1 - \frac{(b + \mu\Delta b)}{\mu(1 - \mu)\Delta b} (\hat{\mu} - \mu) \right) (\mu - \hat{\mu})\Delta b \frac{c}{b + \Delta b \hat{\mu}} \right) = \lambda_{IC} \frac{6c\Delta b(b + \mu\Delta b)(b^2 + \mu(2b + \Delta b))}{\mu(1 - \mu)(b + \Delta b \hat{\mu})^4} > 0$$

and therefore it remains the case that  $\tilde{\ell}''' > 0$  wherever it is continuously differentiable. The kink is concave, as

$$\frac{\partial}{\partial \hat{\mu}} \left( -\lambda_{IC} \left( 1 - \frac{(b + \mu\Delta b)}{\mu(1 - \mu)\Delta b} (\hat{\mu} - \mu) \right) (\mu - \hat{\mu})\Delta b \frac{c}{b + \Delta b \hat{\mu}} \right) \Big|_{\hat{\mu}=\mu} = \frac{c\Delta b \lambda_{IC}}{b + \Delta b \mu} > 0$$

and the curvature of the function increases, as

$$\frac{\partial^2}{\partial \hat{\mu}^2} \left( -\lambda_{IC} \left( 1 - \frac{(b + \mu\Delta b)}{\mu(1 - \mu)\Delta b} (\hat{\mu} - \mu) \right) (\mu - \hat{\mu})\Delta b \frac{c}{b + \Delta b \hat{\mu}} \right) \Big|_{\hat{\mu}=\mu} = -2c\lambda_{IC} \left( \frac{1}{1 - \mu} + \frac{1}{\mu} + \frac{\Delta b^2}{(b + \Delta b \hat{\mu})^2} \right) < 0$$

(Note that the signs are flipped relative to their intuitive interpretation, since the component is part of the Lagrangian for  $\hat{\mu} \leq \mu$ .) This establishes the result, as the pasted Lagrangian is concave

to convex, with a concave kink at  $\mu_i = \mu$ . If  $\tilde{\ell}$  is concave at  $\mu_i = \mu$  with a single support point of concavification, the optimal information structure is uninformative and therefore the IC cannot be satisfied. Therefore, from the shape of  $\tilde{\ell}$ , the concavification is supported at  $\bar{\mu}$  and at a point  $\mu^* < \mu$ , possibly in addition to  $\mu$ .

### A.3 Commitment to a Continuation Value

Suppose that  $u(w) = 2\sqrt{w}$  and that in the first period, the principal can commit to a continuation value, i.e.  $U(s)$ . Note that this does not change the transformation to belief space and we can hence write  $w(\hat{\mu})$ ,  $U(\hat{\mu})$ . The problem reads

$$\Pi_1(\mu) = \max_{m,w} P_\mu Y + \int m(\hat{\mu}) (\delta \Pi_2(\hat{\mu}, U(\hat{\mu})) - w(\hat{\mu})) d\hat{\mu} \quad (46)$$

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

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

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

Straightforward computations establishes that

$$\Pi_2(\hat{\mu}, U(\hat{\mu})) = (a + b + (\Delta a + \Delta b) \hat{\mu}) Y - \frac{U(\hat{\mu})^2}{2} - c^2 \frac{(1 - a - b - (\Delta a + \Delta b) \hat{\mu}) (a + b + (\Delta a + \Delta b) \hat{\mu})}{2 (b + \hat{\mu} \Delta b)^2}.$$

That is, the continuation value is additively separable in the posterior independent cost of providing the continuation value and the cost of providing incentives. This is a feature of the utility function and greatly simplifies the analysis.

Rewriting the contracting problem in utility space, we see that the first period objective reads

$$P_\mu Y + \int m(\hat{\mu}) \left( \delta \left( (P_{\hat{\mu}} Y) - \frac{U(\hat{\mu})^2}{2} - c^2 \frac{(1 - a - b - (\Delta a + \Delta b) \hat{\mu}) (a + b + (\Delta a + \Delta b) \hat{\mu})}{2 (b + \hat{\mu} \Delta b)^2} \right) - \frac{u(\hat{\mu})^2}{2} \right) d\hat{\mu}$$

Equating marginal costs of providing utility to the agent, the optimal contract satisfies  $\delta U(\hat{\mu}) = u(\hat{\mu})$ . Hence, the problem is equivalent to the period by period contracting problem with a cost of utility of  $w(u) = (\delta + \delta^2) \frac{u^2}{2}$ , or, equivalently, a utility function  $u(w) = \frac{2}{(\delta + \delta^2)} \sqrt{w}$ . Therefore, Theorem 1 applies and we have the desired result.

*Remark.* With a general utility function, the costs of providing the continuation utility and the posterior belief are not separable in the principal's continuation profit. This introduces cross-terms in the derivatives of the Lagrangian which are hard to control without making restrictions on the Lagrange multipliers. Hence, the proof strategy of Theorem 1 does not easily generalize to this case.

## 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. 6
- Alonso, Ricardo and Odilon Câmara**, “Bayesian Persuasion with Heterogeneous Priors,” *Journal of Economic Theory*, September 2016, 165, 672–706. 6, 20

- Aumann, Robert J. and Michael Maschler**, *Repeated Games with Incomplete Information*, Cambridge, Mass: MIT Press, 1995. 14
- Bénabou, Roland and Jean Tirole**, “Belief in a Just World and Redistributive Politics,” *The Quarterly Journal of Economics*, May 2006, 121 (2), 699–746. 19
- Bergemann, Dirk and Stephen Morris**, “Information Design: A Unified Perspective,” *Journal of Economic Literature*, March 2019, 57 (1), 44–95. 6
- Bhaskar, V**, “The Ratchet Effect: A Learning Perspective,” Technical Report 2021. 6
- Bhaskar, V. and George J. Mailath**, “The Curse of Long Horizons,” *Journal of Mathematical Economics*, May 2019, 82, 74–89. 6, 25
- Boleslavsky, Raphael and Kyungmin Kim**, “Bayesian Persuasion and Moral Hazard,” *SSRN Electronic Journal*, 2017. 6, 10, 14
- Boretz, Elizabeth**, “Grade Inflation and the Myth of Student Consumerism,” *College Teaching*, 2004, 52 (2), 42–46. 19
- Brunnermeier, Markus K. and Martin Oehmke**, “The Maturity Rat Race,” *The Journal of Finance*, 2013, 68 (2), 483–521. 24
- 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. 4, 19
- Charness, Gary and Uri Gneezy**, “Portfolio Choice and Risk Attitudes: An Experiment,” *Economic Inquiry*, 2010, 48 (1), 133–146. 19
- Cho, In-Koo and David M. Kreps**, “Signaling Games and Stable Equilibria,” *The Quarterly Journal of Economics*, May 1987, 102 (2), 179–221. 24
- **and Joel Sobel**, “Strategic Stability and Uniqueness in Signaling Games,” *Journal of Economic Theory*, April 1990, 50 (2), 381–413. 24
- Dai, Liang, Yenan Wang, and Ming Yang**, “Dynamic Contracting with Flexible Monitoring,” Technical Report 2022. 5
- Datar, Srikant, Susan Cohen Kulp, and Richard A. Lambert**, “Balancing Performance Measures,” *Journal of Accounting Research*, June 2001, 39 (1), 75–92. 5
- de la Rosa, Leonidas Enrique**, “Overconfidence and Moral Hazard,” *Games and Economic Behavior*, November 2011, 73 (2), 429–451. 6, 19
- 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. 6, 25

- Demougin, Dominique and Claude Fluet**, “Monitoring versus Incentives,” *European Economic Review*, October 2001, 45 (9), 1741–1764. 5
- 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. 6
- Doval, Laura and Vasiliki Skreta**, “Constrained Information Design: Toolkit,” Technical Report 2018. 6, 13
- Dumav, Martin, Urmee Khan, and Luca Rigotti**, “Moral Hazard with Heterogeneous Beliefs,” October 2021. 19
- Dye, Ronald A.**, “Optimal Monitoring Policies in Agencies,” *The RAND Journal of Economics*, 1986, 17 (3), 339–350. 5
- 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. 4
- Ekmekci, Mehmet and Nenad Kos**, “Signaling Covertly Acquired Information,” Technical Report June 2021. 25
- Ely, Jeffrey and Martin Szydlowski**, “Moving the Goalposts,” *Journal of Political Economy*, May 2019, p. 704387. 5
- Fang, Hanming and Giuseppe Moscarini**, “Morale Hazard,” *Journal of Monetary Economics*, May 2005, 52 (4), 749–777. 5, 6
- 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. 5
- Filippin, Antonio and Paolo Crosetto**, “Click’n’Roll: No Evidence of Illusion of Control,” *De Economist*, September 2016, 164 (3), 281–295. 19
- Fuchs, William**, “Contracting with Repeated Moral Hazard and Private Evaluations,” *American Economic Review*, September 2007, 97 (4), 1432–1448. 4
- Galperti, Simone**, “Persuasion: The Art of Changing Worldviews,” *American Economic Review*, March 2019, 109 (3), 996–1031. 20
- Georgiadis, George and Balazs Szentcs**, “Optimal Monitoring Design,” *Econometrica*, 2020, p. 55. 5, 6, 14

- 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. [6](#)
- 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. [1](#)
- 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](#), [4](#), [11](#), [30](#)
- Hart, Oliver and Jean Tirole**, “Vertical Integration and Market Foreclosure,” *Brookings Papers on Economic Activity. Microeconomics*, 1990, *1990*, 205. [24](#)
- Hoffmann, Florian, Roman Inderst, and Marcus Opp**, “Only Time Will Tell: A Theory of Deferred Compensation,” *The Review of Economic Studies*, May 2021, *88* (3), 1253–1278. [5](#)
- Holmström, Bengt**, “Moral Hazard and Observability,” *The Bell Journal of Economics*, April 1979, *10* (1), 74–91. [2](#), [4](#)
- , “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](#), [4](#)
- Hörner, Johannes and Nicolas S. Lambert**, “Motivational Ratings,” *Review of Economic Studies*, July 2021. [2](#), [6](#)
- Huffman, David B., Collin Raymond, and Julia Shvets**, “Persistent Overconfidence and Biased Memory: Evidence from Managers,” Technical Report 2019. [4](#), [19](#)
- 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. [19](#)
- 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. [6](#), [14](#)

- Kim, Son Ku**, “Efficiency of an Information System in an Agency Model,” *Econometrica*, 1995, *63* (1), 89–102. [4](#)
- Langer, Ellen J.**, “The Illusion of Control,” *Journal of Personality and Social Psychology*, 1975, *32* (2), 311–328. [4](#), [19](#)
- Larwood, Laurie and William Whittaker**, “Managerial Myopia: Self-serving Biases in Organizational Planning,” *Journal of Applied Psychology*, April 1977, *62* (2), 194–198. [4](#), [19](#)
- Lerner, Melvin J.**, *The Belief in a Just World: A Fundamental Delusion*, Boston, MA: Springer US : Imprint : Springer, 1980. [4](#), [19](#)
- Li, Anqi and Ming Yang**, “Optimal Incentive Contract with Endogenous Monitoring Technology,” *Theoretical Economics*, 2020, *15*, 1135–1173. [5](#)
- Lizzeri, Alessandro, Margaret A. Meyer, and Nicola Persico**, “The Incentive Effects of Interim Performance Evaluations,” Technical Report, Penn Economics Department September 2002. [5](#)
- Luenberger, David G.**, *Optimization by Vector Space Methods* Series in Decision and Control, New York: J. Wiley, 1969. [31](#)
- MacLeod, W. Bentley**, “Optimal Contracting with Subjective Evaluation,” *The American Economic Review*, 2003, *93* (1), 216–240. [4](#)
- 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. [19](#)
- Orlov, Dmitry**, “Frequent Monitoring in Dynamic Contracts,” *Journal of Economic Theory*, forthcoming, p. 54. [5](#), [8](#)
- , **Andrzej Skrzypacz, and Pavel Zryumov**, “Persuading the Principal To Wait,” *Journal of Political Economy*, July 2020, *128* (7). [24](#)
- Park, Young Joon and Luís Santos-Pinto**, “Overconfidence in Tournaments: Evidence from the Field,” *Theory and Decision*, July 2010, *69* (1), 143–166. [19](#)
- Prat, Julien and Boyan Jovanovic**, “Dynamic Contracts When the Agent’s Quality Is Unknown: Dynamic Contracts,” *Theoretical Economics*, September 2014, *9* (3), 865–914. [6](#), [25](#)
- Rodina, David**, “Information Design and Career Concerns,” Technical Report 220 November 2018. [6](#)

**Singer, Natasha**, “In a Mood? Call Center Agents Can Tell,” *The New York Times*, October 2013. 2

**Smolin, Alex**, “Dynamic Evaluation Design,” *American Economic Journal: Microeconomics*, 2021, 13 (4), 300–331. 5

**Treust, Maël Le and Tristan Tomala**, “Persuasion with Limited Communication Capacity,” *Journal of Economic Theory*, November 2019, 184, 104940. 6, 13, 14

**Yaouanq, Yves Le and Peter Schwardmann**, “Learning about One’s Self,” *Journal of the European Economic Association*, forthcoming, p. 39. 4