



Working Paper Series



This paper can be downloaded without charge from:
<http://www.richmondfed.org/publications/>



THE FEDERAL RESERVE BANK OF RICHMOND
RICHMOND ■ BALTIMORE ■ CHARLOTTE

Deferred pay for bank employees: implications of hidden actions with persistent effects in time

Arantxa Jarque* and Edward Simpson Prescott^{†‡}
Federal Reserve Bank of Richmond

Working Paper No. 10-16R

September 10, 2015

Abstract

We develop a two-period model of incentive provision based on Hopenhayn and Jarque (2010). We study the deferral of pay in the optimal contract that implements prudent risk-taking and discuss implications for the regulation of real life compensation instruments, such as bonuses or clawbacks, in the banking sector. In this model, the effects of hidden actions are persistent and hence are revealed over time. We show that, due to the positive effect of persistence on the information available to the principal, a higher intensity of persistence (more influence of the effort on the second period results) will sometimes translate into less deferral of pay (less expected pay in the second period than in the first). We also compare this model with a standard repeated moral hazard, in which the effects of effort do not persist over time. We find that persistence does not necessarily imply a higher proportion of deferred pay.

Key Words: Moral Hazard, Persistence, Bonus, Deferred Pay, Banking.

Journal of Economic Literature Classification Numbers: D80, D82, D86, G21, G28.

**Email:* Arantxa.Jarque@rich.frb.org (Corresponding author)

[†]*Email:* Edward.Prescott@rich.frb.org

[‡]The views expressed in this paper are those of the authors and not necessarily those of the Federal Reserve Bank of Richmond or the Federal Reserve System. This paper is a substantial rewrite of the original WP No. 10-16 in the Richmond Fed working paper series, available here. We would like to thank for their comments Yongsung Chang, Doug Diamond, Huberto Ennis, Joe Haubrich, Boyan Jovanovic, Guillermo Ordoñez, Esteban Rossi-Hansberg, and Nico Trachter. We also benefited from presenting our work at the Stockman Alumni Conference at U. of Rochester, the Midwest Macro Meetings, the SED Meetings, the UVa-Richmond Fed Jamboree, and the Cleveland Fed.

1 Introduction

This paper develops a dynamic model of employee compensation to study the optimal deferral of pay. It builds on the moral hazard model with persistence of Hopenhayn and Jarque (2010) in which actions taken today affect output in the future. We study how deferral varies with the degree of persistence, and how it depends on discounting. We also compare the amount of deferral with that in the benchmark repeated moral hazard model, which has no persistence.

Many actions taken by an employee have persistent effects on an organization. The strategic plan set by a chief executive officer (CEO) affects firm productivity today and in the future. The control environment developed by the chief financial officer (CFO) affects the long-term quality of accounting numbers. The investment made in a client relationship affects sales today and in the future.

These persistent effects are particularly common in banking firms. A loan originated by a mortgage loan officer need not default right away. A credit line originated by a commercial bank officer will generate profits or losses over time. Securities bought by a trader will affect profits until they mature, default, or are sold.

In the recent financial crisis, many loans and investments that were initially profitable became dramatically unprofitable later. There is a wide belief that a reason for this is that bank employees' compensation was separated from the long-term consequences of their actions (FCIC, 2011). As a result, recent regulations in the United States and Europe have focused on regulating banker compensation to control risk. In the United States, much of this regulation has pushed banks to compensate employees through deferring payment of bonuses and allowing bonuses to be clawed back by the firm, with the hope of linking payments to long-term performance (Federal Reserve, 2011).

Our goal is to provide a framework that can be used to assess the desirability of this type of regulations. We develop a two-period moral hazard model in which the hidden actions of a risk-averse agent have persistent effects on output, so information about past actions is revealed over time. In the optimal contract, the principal spreads contingent payments over the two periods to provide incentives. To summarize the time-structure of compensation, we calculate the proportion of the total expected payment over the two periods that will be delivered in the second period. We use this as our measure of deferral and describe how it depends on the intensity of persistence and the discount factor.¹

Our analysis builds on results from Rogerson (1985) and Hopenhayn and Jarque (2010). Rogerson (1985) studied a standard two-period repeated moral hazard problem without

¹A related question of interest is whether an outside observer, such as a regulator, could evaluate compensation plans based on measure such as deferral or the importance of bonuses in overall pay and determine whether they implement "excessive" risk-taking – where a depository institution distorts its decisions due to the safety net. We tackle this question in Jarque and Prescott (2013). See also Phelan (2009). Here, instead, we characterize the optimal contract that implements *prudent* risk-taking (i.e., high effort).

persistence. He showed how, in the presence of commitment to long-term contracts, the optimal smoothing of incentives over time determines the allocation of payments throughout the contract. In a model with only one initial effort, the effects of which persist into the second period, Hopenhayn and Jarque (2010) showed that the optimal contract heavily uses the information generated by returns realized in later periods because this provides efficient incentives for early persistent efforts. This provides a second force that determines the optimal deferral of payments, together with the incentive-smoothing motive characterized in Rogerson (1985), which still prevails in the setting with persistence.

We build on these results by developing a model where the degree of persistence in production and the effectiveness of later actions are indexed by single-dimensional scalars. This way, we can nest two different setups of effort with persistence, as well as the benchmark case of a standard repeated moral hazard model without persistence in Rogerson (1985). We focus our analysis on two particular cases of persistence. In the first one, the agent exerts effort exclusively at the beginning of the contract, but output in the second period is affected by that initial effort. Natural examples in financial firms of such a technology are the jobs of a trader (who exerts effort in selecting only good assets to purchase according to clues about their stochastic future value) and a loan originator (who exerts effort to select borrowers worthy of credit and designs the terms of the loan according to the risk characteristics of the borrowers). This first setup allows us to study the effect on compensation of different degrees of persistence.

The second case we consider allows for a second period action to have an effect on the second period, though the first period effort remains the main determinant in the probability distribution of future output. In financial firms, credit line managers are a good example of such a technology: they exert an initial effort to select borrowers worth of credit and design the terms of the credit line, but they also need to exert effort in future periods to track the evolution of the risk characteristics of the borrowers and adjust the terms of the credit line accordingly. This second setup allows us to study the implications of substitutability of long-term versus short-term effort across periods.

Our results consist of a complete characterization of the contract for the case of an agent with CRRA preferences with a relative risk aversion of one-half. This choice of a functional form allows us to solve for the closed form solution of the contract and to derive the comparative statics. We show that, if only the initial effort affects future output, the optimal amount of deferral may not be monotonically increasing in the intensity of persistence. This is despite the intuitive push for deferral that arises when more information about past efforts is contained in future output. The reason is that there is a second force that reduces deferral. We show that the improved information that comes with increased persistence reduces the overall need for incentives, and that reduces the required amount of deferral. Sometimes the second force outweighs the first, creating the non-monotonicity.

In contrast, in the extension of the persistence model in which future actions also influence

future output, we show that the more important the initial effort is relative to the second, the more deferral there will be in the optimal contract. This is because more persistence now implies a decrease in information about the second period action, making incentives for it more difficult, and hence disproportionately increasing expected consumption in the second period.

Rogerson (1985) showed that, in the standard repeated moral hazard model, deferral was optimal despite the lack of any persistent effects in production. When comparing this model with our main model with persistence, we find that they have different comparative statics of deferral with respect to discounting: more patience implies higher deferral in the repeated moral hazard, while it implies less deferral in the case with persistence. Interestingly, like we found before, the intuition of persistence being the driving force for deferral is incomplete. We find numerous examples in which more compensation is deferred in the repeated moral hazard model than in the persistence model.

In terms of assessing compensation regulations, these results have several implications. These include showing caution when linking persistence to the amount of deferral and that regulations should focus on the actual payouts rather than the legal contractual form. They also demonstrate the need for careful measurement of production technologies.

1.1 Related Literature

The task of studying compensation practices at banking organizations faces an important difficulty in the lack of publicly available data on compensation.² There are a few recent empirical papers about bank employee incentives (Agarwal and Ben-David (2011), Berg, Puri and Rocholl (2014), and Hertzberg, Liberti, and Paravisini, 2011), but they do not address directly the question of deferral of pay. Cole, Kanz, and Klapper (2011) find, in an experimental setting, that a simple form of deferral of incentives (purely postponing payments for three months) makes loan officers exert less screening effort, resulting in lower quality of loans.

On the theory side, a few recent articles have focused on incentives for bank employees. Hartman-Glaser, Piskorski, and Tchisty (2010) focus on the external agency problem between a mortgage originator and the investor buying its securitized mortgages. Inderst and Pfeil (2013) model the internal pay practices of originating firms when the agent needs to both exert a hidden action to generate mortgage opportunities and reveal private information about the quality of those opportunities. They analyze how mandatory deferral of pay

²The SEC requires executive compensation disclosure for the five top-paid executives of publicly traded firms, including financial institutions (Item 402(c)(2)(x) of Regulation S-K). For the rest of the employees of the firms, requirements are minimal. In August of 2015, the SEC adopted a rule, which was mandated by the Dodd-Frank Wall Street Reform and Consumer Protection Act, that requires these same firms to report the ratio of the compensation of their CEO to the median compensation of the employees of the firm starting in 2017.

affects the external agency problem when the firm securitizes the mortgages it originates. Hoffmann, Inderst and Opp (2013) focus in multi-tasking environments in which two hidden actions of the agent (effort and “diligence”) affect the final performance of the firm (which they interpret as mortgage defaults or quality of financial advice). In all of these studies the agent is risk neutral and deferral is costly only because he discounts the future at a higher rate than the principal.

In our paper, we instead assume the agent is risk averse and focus on the case of a common discount factor for the agent and the principal. Hence, our paper relates to the literature on moral hazard problems between a risk-neutral principal and a risk-averse agent when effort has persistent effects in time. In addition to the two models that we directly build on (Hopenhayn and Jarque, 2010; and Mukoyama and Sahin, 2005), related models have been analyzed in Fernandes and Phelan (2000), Kwon (2006), Phelan (2009), and Jarque (2010). These studies use various assumptions on the stochastic processes and the time structure of efforts’ effect over time to handle the technical difficulty that arises with persistence due to the possibility of joint deviations of effort across periods.

1.2 Organization

This article is organized as follows. After describing the environment in Section 2, the optimal contract is characterized in Section 3. Section 4 introduces our measure of deferred pay, and Section 5 presents the results on optimal deferral. After discussing deferral for an extension of the model to two actions in Section 6, Section 7 discusses the implications of the models for the regulation of pay in the financial sector. Section 8 concludes.

2 Environment

There is an agent who operates a two-period project for a principal. In each period, the project produces either a low return, $r_l = 0$, or a high return, $r_h = 1$. The returns are publicly observable and contractible. In general, the agent can take an action, a , in both periods, though we also consider cases in which the second-period action does not affect the outcome of the project. The action is unobservable to the principal, and it can take on two values, a_l or a_h . The agent gets disutility $z(a)$ from taking an action. We set $z(a_l)$ to 0. The disutility of high effort, $z(a_h)$, is set to $\phi > 0$ when effort affects the output and 0 otherwise. When necessary, we abuse notation slightly and use a_1 to refer to the first-period action of the agent and a_2 to refer to the second-period action.

Both the agent and the principal discount the future at a rate $\beta \leq 1$. The principal is risk neutral. The agent is risk averse, with utility of consumption c being $U(c) = \frac{c^{1-\gamma}}{1-\gamma}$, where we set $\gamma = 1/2$ for most of the analysis. Thus, the agent’s utility in each period is $U(c) - z(a)$.

The objective of this paper is to study the deferral of pay. We will assume (for technical reasons) that the agent is not allowed to save or borrow, so pay is equivalent to consumption. It is easy to see that, given the concavity of the utility function, some of the agent’s consumption will be delayed to the second period. We denote as ψ the proportion of expected consumption given in the second period, that is, the proportion deferred. We will study the optimal level of ψ as a function of the primitives.

The actions of the agent affect the return through the stochastic structure of the production technology. We denote the probability of a high return in the first period as $f_1(r_h|a) = \pi(a)$, where:

$$\begin{aligned}\pi(a_h) &= \pi_g, \\ \pi(a_l) &= \pi_b.\end{aligned}$$

We assume that $1 > \pi_g > \pi_b > 0$. The probability of a low output, $f_1(r_l|a)$, is $1 - f_1(r_h|a)$. The probability distribution of a high return in the second period, denoted $f_2(r_h|a_1, a_2)$, depends potentially both on the first and second period actions in the following way:

$$f_2(r_h|a_1, a_2) = \alpha\pi(a_1) + (1 - \alpha)[\sigma\pi(a_2) + (1 - \sigma)\pi_b]. \quad (1)$$

This general structure nests two setups with effort persistence that have been studied in Hopenhayn and Jarque (2010) and Mukoyama and Sahin (2005). It also nests the standard repeated moral hazard model (Rogerson, 1985), which we will use as a benchmark. We will analyze three specific cases defined by different combinations of restrictions on α , σ , and β values.

The stochastic structure of second period output that (1) implies is:

$$\begin{array}{c|c|c} f_2(r_h|a_1, a_2) & a_2 = a_l & a_2 = a_h \\ \hline a_1 = a_l & \pi_b & \alpha\pi_b + (1 - \alpha)[\sigma\pi_g + (1 - \sigma)\pi_b] \\ \hline a_1 = a_h & \alpha\pi_g + (1 - \alpha)\pi_b & \alpha\pi_g + (1 - \alpha)[\sigma\pi_g + (1 - \sigma)\pi_b] \end{array} . \quad (2)$$

Our main case of interest is that of $\sigma = 0$, when the second period effort does not affect output. We focus our analysis on this case, since it provides a clean setup to study the provision of incentives over time, given the timing of revelation of information. In this case, the probability of high output in $t = 2$ is a weighted average of the effort-determined probability and of π_b . The effect of the agent’s effort on the probabilities persists, but it is potentially smaller in $t = 2$ than in $t = 1$. We refer to the exogenous weight of the effort-determined probability, $\alpha \in [0, 1]$, as the “degree” of persistence. We refer to the extreme cases of $\alpha = 0$ and $\alpha = 1$ as no persistence and “perfect persistence,” respectively. We allow for the degree of persistence to vary over the range $\alpha \in (0, 1]$ and study the comparative statics of deferral with respect to this parameter.

The case of $\sigma = 1$ and $\alpha = 0$ is the benchmark repeated moral hazard model (*RMH*) in which only second period effort affects second period output. We analyze this model because, as we will discuss later, it also exhibits deferral of pay despite not having any persistence. By comparing it with the persistence model we will highlight reasons for deferral of pay that are not attributable to persistence.

Finally, we analyze an extension of the main model with persistence in which $\sigma = 1$ and $\alpha > 0$, that is, the second-period action complements the first-period action in determining the probability of the second-period output, and the principal designs the contract to implement $a_2 = a_h$. This constitutes an intermediate case between the pure persistence case of $\sigma = 0$ and the benchmark *RMH*.

3 Optimal contract

We assume that values for r_l, r_h , the cost of effort ϕ , and the stochastic structure are such that the principal finds it profitable to ask the agent to take the high action. Hence, we can just consider the problem of finding the minimal cost to implement a_h . Let $c_i, i = l, h$ be compensation in the first period as a function of the realization of r_i . Similarly, let c_{ij} be compensation in the second period as a function of first period output, r_i , and second period output, r_j . We use the standard relabeling of $u_i = U(c_i)$, and $u_{ij} = U(c_{ij})$, and we denote $h(\cdot) = U^{-1}(\cdot)$. Whenever a second-period action is effective ($\sigma > 0$), the strategy of the agent in the second period may depend on the first period realization. We will denote this strategy as $a_2(r_i)$ for $i = l, h$.

The expected utility of the agent as a function of his strategy of contingent effort choices is:

$$EU(a_1, a_2(r_l), a_2(r_h)) = \sum_i f_1(r_i|a_1) \left[u_i + \beta \sum_j f_2(r_j|a_1, a_2(r_i)) u_{ij} - \beta z(a_2(r_i)) \right] - z(a_1).$$

Recall that whenever second-period effort does not affect output, there is no disutility. Formally:

$$z(a_1) = \begin{cases} \frac{\phi}{1+\beta} & \text{if } \sigma > 0 \\ \phi & \text{if } \sigma = 0. \end{cases} ; \quad z(a_2) = \begin{cases} \frac{\phi}{1+\beta} & \text{if } \sigma > 0 \\ 0 & \text{if } \sigma = 0. \end{cases}$$

This structure of effort disutility helps us to write the problem of the principal in a general form that covers all three cases. In particular, when second-period effort does not affect output ($\sigma = 0$), both the principal and the agent will be indifferent between the low and high actions in the second period. Hence, we can assume they implement a_h in both periods even in this knife-edge case. Note that when $\sigma > 0$ we normalize effort disutility so total effort disutility is the same as in the case $\sigma = 0$.

The principal's optimization problem is:

$$\min_{u_i, u_{ij}} \sum_i f_1(r_i|a_h) \left[h(u_i) + \beta \sum_j h(u_{ij}) f_2(r_j|a_h, a_h) \right] \quad (3)$$

subject to the domain constraints

$$u_i, u_{ij} \geq U(0) \quad \forall i, j, \quad (4)$$

the participation constraint (PC),

$$EU(a_h, a_h, a_h) \geq \bar{U}, \quad (5)$$

where \bar{U} is the reservation utility of the agent, and the incentive constraints (IC).

In the case of $\sigma = 0$ there will simply be one IC, which is for the first-period effort:

$$EU(a_h, a_h, a_h) \geq EU(a_l, a_h, a_h). \quad (6)$$

In the case of $\sigma = 1$ there are six incentive constraints: the two second-period incentive constraints along the equilibrium path

$$\begin{aligned} EU(a_h, a_h, a_h) &\geq EU(a_h, a_l, a_h) \\ EU(a_h, a_h, a_h) &\geq EU(a_h, a_h, a_l), \end{aligned} \quad (7)$$

the first-period incentive constraint (6), and three additional incentive constraints that consider joint deviations of first- and second-period actions:

$$\begin{aligned} EU(a_h, a_h, a_h) &\geq EU(a_l, a_l, a_h) \\ EU(a_h, a_h, a_h) &\geq EU(a_l, a_h, a_l) \\ EU(a_h, a_h, a_h) &\geq EU(a_l, a_l, a_l). \end{aligned} \quad (8)$$

There is a seventh incentive constraint that considers the strategy of (a_h, a_l, a_l) , but it is easy to show that it will be implied by the constraints in (7).

The problem has a strictly convex objective function and linear constraints, so it has a unique solution. In the common solution for the ‘‘first best’’ observable effort case, optimal compensation under full information is characterized by constant consumption over the two periods, since the incentive constraints are not part of the problem. The implied deferred proportion of consumption is $\psi = 1/2$. As we will see when effort is unobservable, compensation may not be evenly split across the two periods due to the trade-off between incentive smoothing motives and the use of later information for incentives. Nevertheless, the property that $\psi > 0$ will be preserved, trivially, whenever the utility of the agent is concave, for consumption smoothing purposes. The focus of our analysis will be characterizing how the degree of persistence of early effort affects the value of ψ .

Next, we analyze the optimal contract when effort is unobservable, for the different parametrizations introduced earlier.

3.1 Benchmark: standard repeated moral hazard without persistence

In our *RMH* model, the agent takes the high action in equilibrium in each of the two periods.³ The probability distribution of return in the second period is, as a special case of matrix 2,

$$\begin{array}{c|cc} f_2(r_h|a_1, a_2) & a_2 = a_l & a_2 = a_h \\ \hline a_1 = a_l & \pi_b & \pi_g \\ \hline a_1 = a_h & \pi_b & \pi_g \end{array}, \quad (9)$$

or, equivalently,

$$f_2(r_h|a_1, a_2) = \pi(a_2).$$

Hence, the probability of high output is π_g if the action is high and π_b if the action is low, i.i.d. across periods.

The first order conditions for this problem are:

$$(u_i) : \frac{1}{U'(c_i)} = \lambda + \mu_1(1 - LR_i) \quad i = l, h$$

$$(u_{ij}) : \frac{1}{U'(c_{ij})} = \lambda + \mu_1(1 - LR_i) + \frac{\mu_{2i}}{\beta f_1(r_h|a_1)}(1 - LR_j) \quad i = l, h, \quad j = l, h,$$

where $\lambda \geq 0$ is the multiplier on the participation constraint, (5), $\mu_1 \geq 0$ is the multiplier of the first-period incentive constraint, and $\mu_{2i} \geq 0$ is the multiplier of the second-period incentive constraint, contingent on the first-period realization being r_i . Let LR_i denote the likelihood ratio of output realization r_i in the first period:

$$LR_i = \frac{f_1(r_i|a_l)}{f_1(r_i|a_h)}, \quad i = l, h.$$

Note that the second-period likelihood ratios, LR_j , are based on individual outcome probabilities, and they are equal to those in the first period:

$$LR_j = \frac{f_2(r_j|a_1, a_2)}{f_2(r_j|a_1, a_2)} = \frac{f_1(r_j|a_l)}{f_1(r_j|a_h)}.$$

It is convenient when presenting our results to denote the variance of the individual outcome likelihood ratios as $v \equiv Var[LR_i|a_h]$, where:

$$v = \frac{(\pi_g - \pi_b)^2}{\pi_g(1 - \pi_g)}.$$

³This is a standard textbook model. For details, see Rogerson (1985).

The solutions for the optimal utilities when $U(c) = \sqrt{c}$ (assuming the lower bound of utility in (4) is not binding) are

$$u_i = \frac{1}{2} \left\{ \frac{\bar{U} + \phi}{1 + \beta} + \frac{\phi}{1 + \beta} \frac{1 - LR_i}{v(1 + \beta)} \right\} \quad (10)$$

$$u_{ij} = \frac{1}{2} \left\{ \frac{\bar{U} + \phi}{1 + \beta} + \frac{\phi}{1 + \beta} \left(\frac{1 - LR_i}{v(1 + \beta)} + \frac{1 - LR_j}{v} \right) \right\}, \quad (11)$$

3.2 Main model: initial effort only

In this case, output is only affected by the first-period action, so $\sigma = 0$. The simplified production technology is:

	$f_1(r_h a_1)$	$f_2(r_h a_1, a_2)$
$a = a_h$	π_g	$\alpha\pi_g + (1 - \alpha)\pi_b$
$a = a_l$	π_b	π_b

Because the second-period probability is only conditional on the first-period action, we introduce a simpler notation:

$$f_2(r_h|a_1, a_2) = g_2(r_h|a_1).$$

As in the previous model, LR_i denotes the likelihood ratio of output realization r_i in the first period. Let LR_{ij} denote the likelihood ratio of the history of first- and second-period output realizations, (r_i, r_j) :

$$LR_{ij} = \frac{f_1(r_i|a_l)g_2(r_j|a_l)}{f_1(r_i|a_h)g_2(r_j|a_h)}, \quad i, j = l, h.$$

The first-order conditions of this problem are

$$\begin{aligned} \frac{1}{U'(c_i)} &= \lambda + \mu(1 - LR_i), \quad \forall i, \\ \frac{1}{U'(c_{ij})} &= \lambda + \mu(1 - LR_{ij}), \quad \forall i, j, \end{aligned} \quad (12)$$

where $\mu \geq 0$ is the multiplier on the incentive constraint, (6). The following proposition follows from these first-order conditions:

Proposition 1 (Hopenhayn and Jarque, 2010) *For any utility specification that satisfies $U' > 0$, $U'' < 0$, and assuming the domain constraint for utility (4) does not bind, compensation levels in the optimal contract are ranked by likelihood ratios:*

$$\begin{aligned} c_h &> c_l, \text{ since } LR_h < LR_l, \\ c_{ij} &> c_{i'j'} \iff LR_{ij} < LR_{i'j'}. \end{aligned}$$

Collorary 1 *In particular,*

If $\alpha = 0$ then $c_{ll} = c_l = c_{lh} < c_{hl} = c_h = c_{hh}$.

If $\alpha \in (0, 1)$ then $c_{ll} < c_l < c_{lh}$, and $c_{hl} < c_h < c_{hh}$.

If $\alpha = 1$ then $c_{ll} < c_l < c_{lh} = c_{hl} < c_h < c_{hh}$.

Note that the solution for consumption does not depend on the period when the transfer takes place except through the likelihood ratios. Another important property of the optimal contract is that it never implies perfect insurance in the first period; information is used as soon as it becomes available.

It will be useful to focus on the case of perfect persistence to discuss the intuition of the results. Here, when $\alpha = 1$, we see that the above proposition implies $c_{lh} = c_{hl}$. Also, incentives are more high powered in the second period after a low return in the first period. That is, it is easy to see that, for any arbitrary functional form of the utility of consumption,

$$0 < u_{hh} - u_{hl} < u_h - u_l < u_{lh} - u_{ll}.$$

For our utility specification of $U(c) = 2\sqrt{c}$, when the lower bound of utility in (4) is not binding we can solve for the closed form solution for the contingent utilities (see Appendix 1 for details of the derivation):⁴

$$u_s = \frac{\bar{U} + \phi}{1 + \beta} + \phi \frac{1}{v_1 + \beta v_2} (1 - LR_s), \quad (13)$$

where $s \in \{l, h, ll, lh, hl, hh\}$, $v_1 \equiv Var [LR_i|a_h]$ denotes the variance of the likelihood ratios in the first period (which is equal to v given our assumptions about the technology in the first period), and $v_2 \equiv Var [LR_{ij}|a_h]$ denotes the variance of the likelihood ratios in the second period, which depends on α .

3.3 Extension: two efforts

In this case, output is affected by the first- and second-period actions, so $\sigma = 1$. Because we want to analyze the effects of persistence on deferral of pay, we limit our analysis of this case to parametrizations in which the relative effectiveness of the first-period effort is high relative to that of the second-period effort. In particular, we restrict $\alpha \in (0.5, 1)$, and $\beta \geq \frac{1-\alpha}{\alpha}$. As shown in Mukoyama and Sahin (2005), this combination of parameters implies a particularly simple form of the optimal contract in spite of the repeated action with persistence.

The principal designs the optimal contract to implement $a_1 = a_h$, and $a_2 = a_h$. The first-period simplified production technology is the same as in the main model. The simplified second-period production technology is:

⁴If the domain constraint in (4) were to bind, some of the strict inequalities in Prop. 1 would become equalities.

$$\begin{array}{c|c|c}
f_2(r_h|a_1, a_2) & a_2 = a_l & a_2 = a_h \\
\hline
a_1 = a_l & \pi_b & \alpha\pi_b + (1 - \alpha)\pi_g \\
\hline
a_1 = a_h & \alpha\pi_g + (1 - \alpha)\pi_b & \pi_g
\end{array} \quad , \quad (14)$$

where α represents the relative importance of the first-period effort in the probability of the high output in period 2.

Recall that the relevant incentive constraints for this extension with two efforts are the first-period incentive constraints in (6) and (8). Fortunately, these six incentive constraints can be reduced to one following the analysis in Mukoyama and Sahin (2005). The authors present conditions that are sufficient for the four first-period deviations to not be profitable for the agent, given that the second-period constraints in (7) are satisfied (see Prop. 4 in their paper). This implies perfect insurance in the first period, which in turn makes the two constraints in (7) become equivalent. In our setup, the marginal effect of a deviation in one period is independent of the choice of action in the other period. This simplifies the analysis further, and it gives rise to the following condition (the counterpart of the condition in Prop. 4 of Mukoyama and Sahin (2005) in our setting).

Lemma 1 *The optimal contract with two efforts and persistence exhibits perfect insurance in the first period if and only if $\beta \geq \frac{1-\alpha}{\alpha}$.*

To characterize the optimal contract when the domain constraint in (4) is not binding we can use the fact that the first-period constraints do not bind and there is perfect insurance in the first period:

$$\begin{aligned}
u_l &= u_h \equiv u_1, \\
u_{lh} &= u_{hh} \equiv u_{2h}, \\
u_{ll} &= u_{hl} \equiv u_{2l}.
\end{aligned}$$

Hence, all variable compensation is deferred until the second period. Intuitively, for the first period deviations to be non-binding we need the loss in expected utility of the agent in the second period to be greater when he chooses to deviate in the first period than when he chooses to deviate in the second one:

$$\beta (u_{2h} - u_{2l}) \alpha (\pi_g - \pi_b) \geq (u_{2h} - u_{2l}) (1 - \alpha) (\pi_g - \pi_b). \quad (15)$$

When this inequality is satisfied, incentives for the first-period action are provided entirely by incentives for the second-period action in the second period. This is the relevant comparison both on the equilibrium path (constraint (6)) and off the equilibrium path (constraint (8)) because the probability is linear in both efforts, and hence the marginal effect of a deviation in the first period is independent of the action choice in the second period. The relation between α and β stated in the lemma follows from (15).

Intuitively, the second-period action is not very effective at increasing the probability of a high output, so implementing the high action in the second period requires high-powered incentives. However, the first-period action is very effective in increasing the probability of high output in the second period, so the agent works hard in the first period to earn a high pay in the second. Incentives in the first period come “for free” in this setup.

The conditions of $\alpha \in (1/2, 1]$ and $\beta \in [\frac{1-\alpha}{\alpha}, 1]$ satisfy the conditions of Lemma 1.⁵ We can now fully solve for the optimal contract. The constraints of the problem simplify to the PC

$$u_1 + \beta [\pi_g u_{2h} + (1 - \pi_g) u_{2l}] - \phi = \bar{U} \quad (16)$$

and the second period ICs, which are both equal and hence simplify to a unique IC:

$$\begin{aligned} \pi_g u_{2h} + (1 - \pi_g) u_{2l} - \frac{\phi}{1 + \beta} &= [\alpha \pi_g + (1 - \alpha) \pi_b] u_{2h} \\ &+ [1 - \alpha \pi_g - (1 - \alpha) \pi_b] u_{2l}. \end{aligned} \quad (17)$$

Because compensation in the second period does not depend on the whole history of output, the relevant likelihood ratio is that of just the second-period output:

$$LR_j = \frac{f(r_j | a_h, a_l)}{f(r_j | a_h, a_h)}.$$

The first-order conditions of the problem are:

$$\begin{aligned} (c_1) &: \frac{1}{U'(c_1)} = \lambda \\ (c_j) &: \frac{1}{U'(c_i)} = \lambda + \frac{\mu}{\beta} [1 - LR_j], \end{aligned}$$

where μ is the multiplier of the IC constraint in Eq. (17). For $U(c) = 2\sqrt{c}$, the constraints together with the first-order conditions imply that

$$\begin{aligned} u_1 &= \frac{\bar{U} + \phi}{(1 + \beta)} \\ u_{2i} &= \frac{\bar{U} + \phi}{(1 + \beta)} + \frac{\phi}{(1 + \beta)} \frac{\beta}{(1 - \alpha)^2} \frac{1}{v} (1 - LR_i), \text{ for } i = l, h. \end{aligned} \quad (18)$$

4 Measuring deferred pay

In this section, we discuss our choice for the measure of deferred pay, ψ , its empirical counterpart, and its suitability as a measure of deferred incentives. We argue that ψ can help us

⁵Mukoyama and Sahin (2005) present a more general probability structure and then proceed to characterize the subset of probability and discount factor combinations that imply perfect insurance in the first period.

understand how persistence of effort and the particular structure of the incentive problem affect — and, in particular, postpone — the optimal provision of incentives over time.

As we stated earlier, the optimal contract with commitment will always trivially imply deferred pay, i.e., the proportion of pay delivered in the second period, ψ , will always be strictly greater than 0. What we want to study is the determination of the optimal level of ψ and its comparative statics with respect to the key parameters of the model. For our analysis, it will be convenient to express deferral as the ratio of second period expected consumption over that of the first. We denote this ratio as D :

$$D \equiv \frac{E[c_2]}{E[c_1]} = \frac{\psi}{(1-\psi)},$$

where the expectation is taken with respect to the effort choices in both periods in equilibrium. Note that D is simply a monotonic transformation of ψ . Furthermore, when $\psi = 0$, $D = 0$, and when $\psi = 1/2$, $D = 1$.

The empirical appeal of this measure is that it can be constructed even by a regulator who only observes realized wages and bonuses, rather than the ex-ante details of incentive contracts. With a large number of bankers, D corresponds to the ratio of average compensation in the second period over average compensation in the first period.

For the particular utility function specification of $U(c) = \sqrt{c}$, the ratio of expected consumptions informs us about the time allocation of variable pay. We can show that a higher D corresponds to more variable pay in the second period by deriving the Rogerson condition from the first order conditions in the previous section. For all three cases (benchmark RMH, main model, and extension), this condition is:

$$\frac{1}{U'(c_i)} = E \left[\frac{1}{U'(c_{ij})} | a_h, c_i \right] \quad \forall i, j. \quad (19)$$

For $U(c) = 2\sqrt{c}$, this becomes

$$E[U(c_i)] = E[U(c_{ij})]. \quad (20)$$

In our analysis of the optimal contracts we will be able to characterize the effects of changes in persistence on the variance of utility in each of the two periods. By the concavity of U and Jensen's inequality, equation (20) implies that an increase in the variance of utility will imply an increase in expected consumption. This will allow us to use our results to determine the allocation of variable pay across periods.⁶ In what follows, we study these comparative statics on our measure D for the three cases of our model.

⁶Under certain implementations of the optimal contract with real-life compensation instruments, variable pay will correspond to “bonus” payments (see section 7 for detailed discussion of this mapping). Hence, our measure allows us to evaluate proposals of deferring “bonuses” or variable pay.

5 The determinants of deferral of pay

We now study the determinants of deferral of pay in our main case of interest, where only the first-period effort affects output in both periods. First, we point out the effect of risk aversion of the agent in the relative importance of expected consumption across the two periods. This is a well-known result since Rogerson (1985), and we simply lay out the implications for our measures of deferred pay given our assumption of $U = 2\sqrt{c}$. Then, we analyze comparative statics of our measure D with respect to the intensity of persistence, α , and then with respect to the discount factor, β .

The degree of risk aversion of the agent, or the concavity of his utility function, will determine whether expected consumption will increase or decrease across the two periods in the optimal contract. This is a direct consequence of the optimal structure of contingent pay determined by the inverse Euler equation in (19), also known as the Rogerson condition, first discussed in Rogerson (1985).

Proposition 2 (Rogerson, 1985) *When the principal and the agent discount the future at the same rate, $E[c_1] < E[c_2]$ whenever $\frac{1}{U'(c)}$ is concave; the inequality is reversed whenever $\frac{1}{U'(c)}$ is convex, and it becomes equality whenever $\frac{1}{U'(c)}$ is linear (for example, with $U(c) = \ln(c)$).*

It follows from this result that, within the class of CRRA utility functions, of the form $\frac{c^{1-\gamma}}{1-\gamma}$, coefficients of relative risk aversion (γ) of one or less imply an increasing path for expected compensation, while coefficients of one or more imply a decreasing path. Hence, for an agent with a relatively low risk aversion, as our $U = 2\sqrt{c}$ implies, our measures of deferral ψ and D will be relatively low.⁷ We will discuss how our comparative statics below may be affected by changes in the degree of risk aversion.

Using the inverse of the utility function, $h(U)$, and the solution for contingent utilities in (13), we get the following expression for our measure of deferred pay:

$$D = \frac{(\bar{U} + \phi)^2 + \left(\frac{\phi}{v_1 + \beta v_2}\right)^2 v_2}{(\bar{U} + \phi)^2 + \left(\frac{\phi}{v_1 + \beta v_2}\right)^2 v_1}. \quad (21)$$

⁷There is a well-known but important implication that stems from equation (19). The optimal contract is such that the agent is always left with a desire to save (assuming that the savings rate is the inverse of the discount rate, $1/\beta$), even if the path for expected compensation is increasing, as it is in our example with $U(c) = 2\sqrt{c}$. Were the agent allowed to save, the actual consumptions would be more back loaded, affecting our measure of deferral. Solving for the path of contingent consumptions in this case, however, is not easy. The incentive constraints are harder to analyze because they need to take into account both effort and savings deviations. Most of the literature on moral hazard and hidden savings (Abraham and Pavoni, 2006; Kocherlakota, 2004), assume a continuum of efforts and look for conditions under which the first order condition of the maximization problem of the agent identify a global maximum. In our setting with binary effort choice that strategy is not valid. Hence, we assume the agent is not able to save.

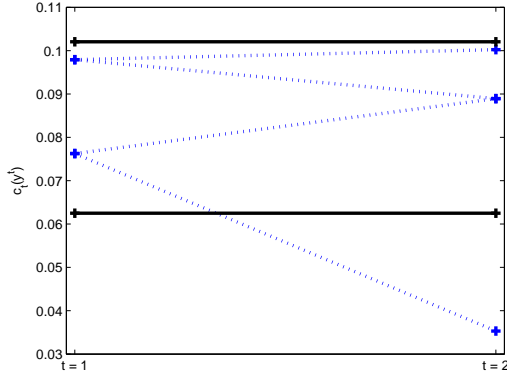


Figure 1: Contingent time paths of consumption in the optimal contract for the main model (solid line correspond to $\alpha = 0$, dotted line corresponds to $\alpha = 1$). Parameters: $\gamma = 0.5$, $\bar{U} = 1$, $\phi = 0.1$, $\pi_g = 0.8$, $\pi_b = 0.4$.

Given the time structure of the problem, the information revealed in the first period, which is captured in the variance of the likelihood ratios in the first period, v_1 , is always available to the principal in the second period. Hence, trivially, we always have $v_2 \geq v_1$. As the Lemma 2 states, a second-period realization that is affected by the initial effort of the agent will always strictly increase the information available in the second period.⁸

Lemma 2 *We have $v_2 > v_1$ if and only if $\alpha > 0$.*

Because utility is linear in the likelihood ratios (see the first order conditions in (12)), it is easy to show that the variance of utility will also be higher in the second period.

Figure 1 illustrates the typical structure of contingent payments in the optimal contract, and how they are affected by persistence. It compares the optimal contract for the case of no persistence ($\alpha = 0$, solid lines) with that of perfect persistence ($\alpha = 1$, dotted lines). Consider first the benchmark case of no persistence ($\alpha = 0$), in which the second-period observation does not depend on the effort of the agent, but rather it is purely random. Under this parametrization, the optimal contract consists of the same contingent consumptions in the two periods. This implies a deferral measure of $D = 1$, or, equivalently, $\psi = 1/2$.

In contrast, consider now the case with perfect persistence ($\alpha = 1$), in which the second period observation does depend on the effort of the agent. The optimal contract now

⁸The results in this section relate to the literature on the “quality of information” in moral hazard models. See Grossman and Hart (1986) for a seminal contribution that ranks distributions in terms of “garbling,” and Kim (1995) for a refinement that proposes mean preserving spreads rankings of the distribution of the likelihood ratios

prescribes different consumptions in the second period, for the same period realization, depending on what the first-period realization was. Looking at first-period consumption in Figure 1 we can see that, in the first period, the “reward” for a high output realization is not as large as it was in the case of $\alpha = 0$, and the “punishment” for a low one is not as severe. On the other hand, comparing second-period consumption after a low realization in the first period, the punishment for another low output is very severe. Next, we analyze how this time structure of payments due to persistence translates into deferral.

Higher persistence implies higher variance of the likelihood ratios in the second period and higher variance of utility. However, this does not necessarily translate into more deferred pay, as measured by D , as the next proposition states.

Proposition 3 *The amount of pay deferred may increase or decrease when the intensity of persistence increases, i.e., $\frac{\partial D}{\partial \alpha} \gtrless 0$.*

The reason for the non-monotonicity is best understood by looking at the expression for D in terms of the Lagrange multipliers of the cost minimization problem:

$$D(\alpha) = \frac{\lambda^2 + [\mu(v_1, v_2(\alpha))]^2 v_2(\alpha)}{\lambda^2 + [\mu(v_1, v_2(\alpha))]^2 v_1}$$

The proof of the proposition shows that the sign of the derivative $\frac{\partial D}{\partial \alpha}$ is positive if

$$\mu(\lambda^2 + \mu^2 v_1) \frac{\partial v_2}{\partial \alpha} + 2\lambda^2 \frac{\partial \mu}{\partial \alpha} (v_2 - v_1) > 0. \quad (22)$$

This expression shows that an increase in persistence has two opposite effects on deferral. The first effect is a force for backloading compensation. Since, as the proof shows, $\frac{\partial v_2}{\partial \alpha} > 0$ due to more information content in the second-period output, the first term in (22) is positive. The second effect is a force for frontloading compensation. Since, as the proof shows, $\frac{\partial \mu}{\partial \alpha} < 0$, and $v_2 > v_1$ by Lemma 2, the second term in (22) is negative. Intuitively, an increase in α , which increases v_2 , implies that providing incentives for high effort will now be achieved more efficiently, reducing the multiplier on the incentive constraint. The term $\frac{\partial \mu}{\partial \alpha} (v_2 - v_1)$ being negative reflects the fact that, since expected utility is equated across the two periods, and the variance of utility is higher in period two, a decrease in the need for incentives in the contract implies a larger reduction in expected consumption in the period where utility is more variable. The terms that multiply the main derivatives of interest mainly capture the effect of the level of expected consumption ($E[c_1]$) in the race between the two forces.

When the second effect is stronger than the first, deferral decreases with the degree of persistence. Figure 2 presents the comparative statics for Example 1, a case where the amount of deferral is non-monotonic in α .⁹ All our numerical examples suggest the same

⁹Parameter values for Example 1: $\bar{U} = 0.3$, $\phi = 0.1$, $\pi_g = 0.8$, $\pi_b = 0.2$, $\beta = 0.95$, and $U(c) = \frac{c^{1-\gamma}}{1-\gamma}$ with $\gamma = 0.5$.

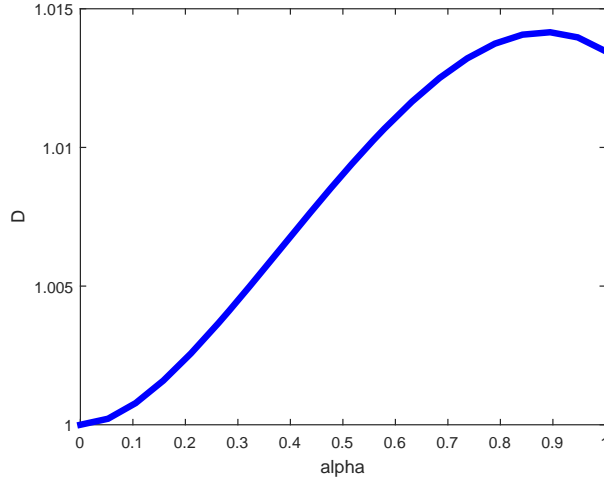


Figure 2: Comparative statics for Example 1 showing non-monotonicities in deferral as a function of α , for $U(c) = \frac{c^{1-\gamma}}{1-\gamma}$, with $\gamma = 0.5$.

single-peaked pattern for all non-monotonic cases, i.e., that the derivative turns negative for the higher range of α . In fact, we can show that $\frac{\partial D}{\partial \alpha} > 0$ always at $\alpha = 0$, implying that there is always an increasing portion of deferral (see Lemma 3 in Appendix 2). Figure 3 presents, for Example 1, the behaviour of expected consumption, the variance of consumption, as well as for the variance of the likelihood ratios in period 2, v_2 , and the multiplier of the incentive constraint, μ . We see that, as α increases, v_2 increases and μ decreases, consistent with the results in Prop. 3. Although we cannot show formally that this is the case, we see in the top right panel that the variance of consumption in each period decreases with α . This was the case in all our numerical examples.

As we pointed out earlier, our measure of deferral will be affected by the degree of risk aversion. Although analytical characterizations for other relative risk aversion coefficients are not available, based on numerical examples we conjecture that the non-monotonicity that we just characterized survives to changes in this parameter. Figure 4 presents such an example.¹⁰ For a relative risk aversion of 1.2, as dictated by proposition 2, $D < 1$ whenever $\alpha > 0$. In this case, deferral first decreases with α , then increases. As Figure 5 shows, the forces behind the non-monotonicities are similar in this second example than in our first example for $\gamma = 0.5$: except for the fact that $E[c_1] > E[c_2]$ in example 2, all other variables change qualitatively in a similar way.

Given that the optimal contract will always involve deferring some incentives, but both the principal and the agent discount future payments, it is interesting to study what is the

¹⁰Parameter values for Example 2: $\pi_g = .8$, $\pi_b = .5$, and $U(c) = \frac{c^{1-\gamma}}{1-\gamma}$, with $\gamma = 1.2$.

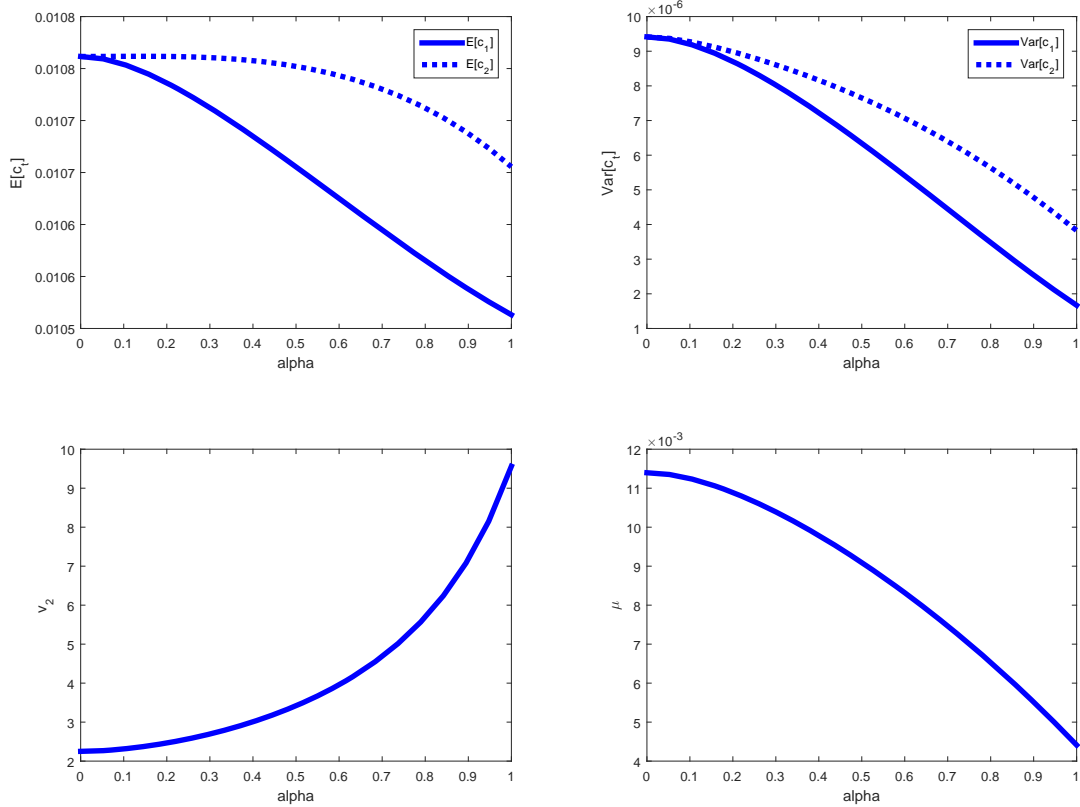


Figure 3: Comparative statics with α for Example 1: (top left panel) expected consumption in each period, (top right panel) variance of consumption in each period, (bottom left panel) variance of the likelihood ratios in the second period, and (bottom right panel) multiplier of the incentive constraint.

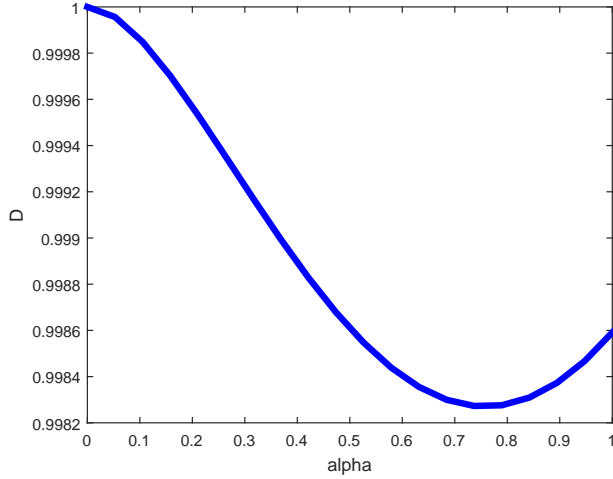


Figure 4: Comparative statics for Example 1 showing non-monotonicities in deferral as a function of α , for $U(c) = \frac{c^{1-\gamma}}{1-\gamma}$, with $\gamma = 1.2$.

effect of both players becoming less impatient.

Proposition 4 *The proportion of pay deferred in the main model decreases with the discount factor β .*

Because of symmetric discounting by the principal and the agent, an increase in β does not make payments in the second period relatively more attractive to the principal. However, it does make the information revealed in the second period more useful: when the agent decides on his effort choice in the first period, he weighs more the effect that it will have in his second period utility. Hence, incentives are cheaper overall because the existing information can be used more effectively. Formally, the value of the IC multiplier, μ , decreases. By a similar logic to the one discussed above, a decrease in the need for incentives in the contract implies a larger reduction in expected consumption in the period where utility is more variable. This decreases expected consumption in the second period more than in the first, decreasing D .

5.1 Comparison to benchmark of RMH

Even without persistence, there will be deferral in the optimal contract for our RMH benchmark. However, the reasons for this deferral are solely to smooth incentives over time, since second-period output provides no information about first-period actions. The dependence of second-period payments on first-period outcomes is optimal because it helps reduce the variance of compensation in the first period, reducing the principal's expected payments.

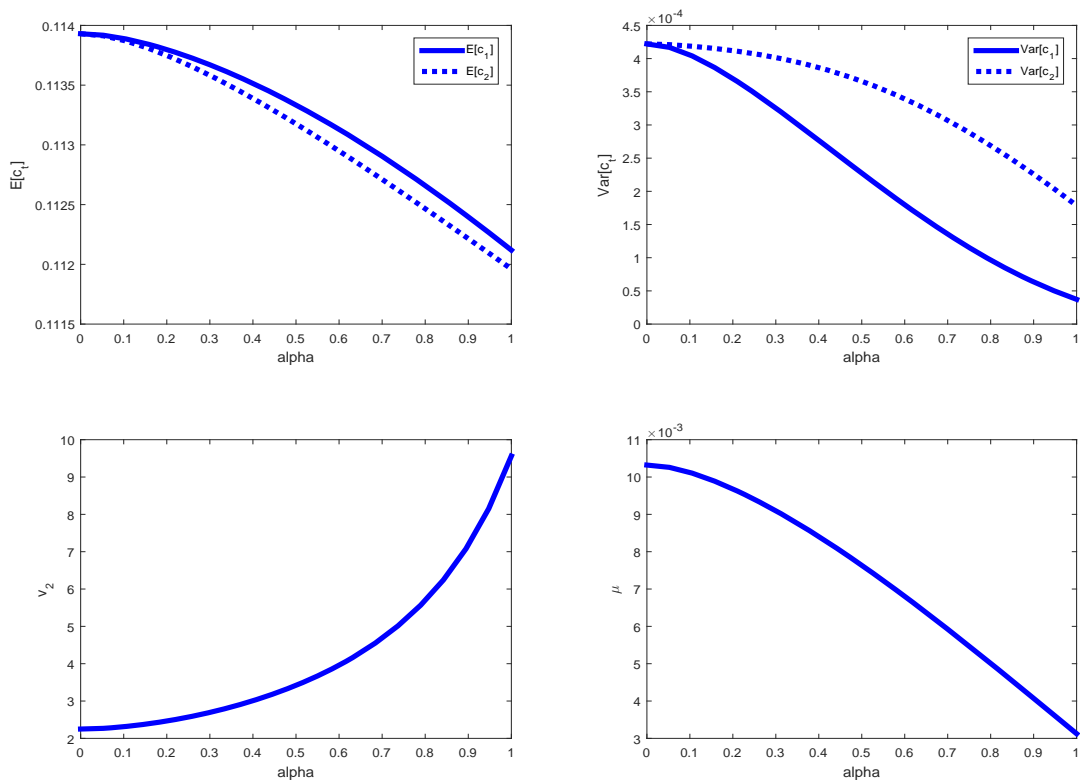


Figure 5: Comparative statics with α for Example 2: (top left panel) expected consumption in each period, (top right panel) variance of consumption in each period, (bottom left panel) variance of the likelihood ratios in the second period, and (bottom right panel) multiplier of the incentive constraint.

We show in this section that this implies that the comparative statics with respect to the discount factor will differ from those in the main model.

With some algebra that uses the closed form solutions in equations (10) and (11), we can show that the portion of the second-period expected consumption that is attributable to incentives for the first-period effort can easily be isolated. This is due to our utility function specification, but it gives us insight into the comparative statics of deferred pay in this model. Deferred pay for the *RMH* model is:

$$D^{RMH} = \frac{E[c_1] + \frac{z(a_2)^2}{v}}{E[c_1]}, \quad (23)$$

or, in terms of the primitives, is

$$D^{RMH} = 1 + \frac{\frac{z(a_2)^2}{v}}{\left[\frac{\bar{U}}{(1+\beta)} + z(a_1) \right]^2 + \frac{z(a_1)^2}{v} \frac{1}{(1+\beta)^2}}.$$

(Note that, given that the principal implements $a_1 = a_2 = a_h$, we will have $z(a_1) = z(a_2)$.)

Proposition 5 *The proportion of pay deferred in the benchmark RMH model increases with the discount factor β .*

In the *RMH* model, a change in the discount factor does not affect the spread $u_{iH} - u_{iL}$, for $i = L, H$. In other words, an increase in the discount factor does not affect incentives at all in the second period in the *RMH* model. It does, however, make first-period incentives cheaper because of a greater ability to spread them over the two periods, implying $\frac{\partial D}{\partial \beta} > 0$. This is in contrast with the effect of β in the main model, where $\frac{\partial D}{\partial \beta} < 0$. In that model, we saw that an increase in the discount factor affected the cost of incentives in all periods proportionally, because incentives are for one effort only and those are spread over the two periods according to the relative amount of information in each.

Given the above result, it is not surprising that we can find examples for which deferral will be higher in the *RMH* model than in the model with persistence. One such example is:

Example 3 Whenever $\bar{U} = 0$, $\alpha = 1$, $z(a_1) = z(a_2) = \phi / (1 + \beta)$, and $\beta = 1$, the proportion of pay deferred in the *RMH* model is always higher than in the main model. (See Appendix 2 for details.)

6 Extension: Deferral with two actions and persistence

We now turn to studying the comparative statics of deferral in our extension with two actions. This exercise highlights that implications for deferral will depend on the particular structure of asymmetric information. We find that, in this particular case of persistence with two actions, there will never be a decrease in deferral when α increases.

Proposition 6 *In the model with repeated actions and strong persistence ($\alpha \in (.5, 1)$ and $\frac{1-\alpha}{\alpha} \leq \beta$), the amount of pay deferred always increases when the intensity of persistence increases, i.e., $\frac{\partial D}{\partial \alpha} > 0$.*

A sketch of the proof, which is included in Appendix 2, is as follows. Deferral can be written as

$$D = \frac{\lambda^2 + [\mu(\alpha)]^2 \tilde{v}_2(\alpha)}{\lambda^2},$$

or, in terms of v , as

$$D = \frac{\lambda^2 + [\mu(\alpha)]^2 (1 - \alpha)^2 v}{\lambda^2}.$$

Our characterization in Section 3.3 showed that pay is only variable in the second period. It is easy to see that the variance of the likelihood ratios, $\tilde{v}_2 = (1 - \alpha)^2 v$, decreases with α , lowering expected consumption in the second period and lowering deferral. However, when we increase α , even if we are increasing persistence, we also decrease the effect of the second-period effort on output. This makes incentives for the second-period effort more costly to implement, that is, μ increases when α increases, raising deferral. We can show that this second effect is stronger, making deferral increase with α in this setup.¹¹

7 Implications for the regulation of deferred pay

Many of the recent regulations on banker compensation focus on the “bonuses” paid to bank employees, with requirements that portions of bonuses be deferred and be clawbacked in the future if long-term performance suffers. The idea behind the regulations is that by better tying bonus payments to long-term performance, bank risk will be reduced (in our model, this corresponds to implementing the high effort because that raises the probability that the loan is repaid). What our analysis underscores is that it is not the legal form of payments to workers that matters for incentives, but the actual state-contingent payments. One implication is that the optimal deferral may very well be in place even in the absence of formal escrow accounts with clawback provisions for bonuses. In other words, because there is a wide variety of legal compensation arrangements that implement a collection of state-contingent payments, regulations that focus on one particular arrangement may be ineffective.

To illustrate this point, we provide two different compensation arrangements that implement the optimal contract in the persistence model: one that uses deferred bonuses and clawbacks and one that does not.¹²

¹¹For values of α outside of the restrictions stated in Proposition 6, for which the perfect insurance result does not hold, we ran a wide range of numerical examples, and in all of these we found $\frac{\partial D}{\partial \alpha} > 0$.

¹²For more discussion along these lines and other examples of implementation see a previous version of this working paper, available here.

The first compensation scheme we propose (compensation package 1) uses a combination of a wage, first-period bonuses, deferrals, and clawbacks, much like the legal form of the contracts considered by the regulations. Let W be the constant base wage. Let B_i be the contingent bonus payment, let d_i be the contingent fraction of the bonus that is deferred, and let θ_i be the contingent fraction of the bonus that is clawed back if second-period output is low; the subscript $i = l, h$ indicates that all these instruments are contingent on the first-period output. In summary, output-contingent compensation under this arrangement after realization r_i in the first period is $W + (1 - d_i) B_i$. Following realization r_h in the second period, it is $W + d_i B_i$, while following r_l it is only $W + d_i B_i (1 - \theta_i)$. That is, if the second-period realization is low, a portion θ_i (where i denotes the first-period realization of r) is clawed back and the agent cannot consume it.

One can show that, when using compensation packages of the type described in scheme 1, perfect persistence ($\alpha = 1$) calls for higher deferral than non-persistence ($\alpha = 0$), and persistence calls for clawbacks, while non-persistence does not. (Section 9.4 in the Appendix contains the details of this implementation and the proof of this claim). Based on this example, one could conclude that a simple way to regulate compensation in occupations that have the feature of information being revealed over time is to demand a high level of the amount of deferral, d , and to use clawbacks. However, there exist alternative implementations of the optimal contingent scheme in which such mandatory high levels of deferral would be distorting. One example would be the arrangement we propose next, compensation package 2.

In compensation package 2, second-period wages are dependent on first-period outcomes, that is, we allow for W_i , $i = l, h$. These contingent wages in the second period can be interpreted as promotions or demotions as a function of early performance. A first-period bonus, B_0 , is paid in addition to W in the first period only if output is high. Following output r_i in the first period, contingent payments in the second period are W_i if output is low and $W_i + B_0$ if output is high. The derivations in Section 9.4 in the Appendix show that, by having $W_l < W < W_h$ and a first period bonus scheme independent of that of the second period, there is no need to have explicit deferral or clawback of bonuses. Even in the case of perfect persistence, $\alpha = 1$, just because there are no formal clawbacks in place one cannot conclude that incentives are not properly implemented

Needless to say, there are a multitude of other ways to implement any desired state-contingent compensation levels. There are guaranteed bonuses, varying performance targets for bonus plans, changes to non-pecuniary benefits, etc. There is also ample anecdotal evidence that bonuses and raises depend on subjective evaluations of the worker's performance. This implies that managers can make compensation implicitly depend on past performance and in ways that make it hard for an outsider to evaluate.

Our discussion in this section calls for caution when evaluating and regulating compensation schemes. First, in terms of evaluation, the formal language used to describe compen-

sation, such as “bonus,” does not describe the actual state–contingent payments, which are what matter for incentives. Second, in terms of regulation, rules that impose clawbacks or reduce bonus payments may be ill–suited, given the ability of firms to use other instruments to implement state–contingent compensation levels. Our discussion does say that evaluation and regulation should be based on measurement of actual payouts, and measures like the D measure we proposed should be useful along this dimension. Nevertheless, it should also be recalled that for such a measure of deferral we have shown that contract theory does not necessarily call for more deferral when the intensity of persistence is higher. In fact, we have shown that a good understanding of the particular structure of the incentive problem is necessary to understand the right amount of deferral that should be expected in the prudent risk–taking contract. In particular, measurement of the amount of deferral is not the only thing required, but detailed measurement of the details of the production technology, e.g., how much persistence is there, is also necessary. A one–size–fits–all approach is likely not appropriate to regulate deferral of pay in the financial industry.¹³

8 Conclusion

This paper characterizes the optimal deferral of compensation in a two–period principal–agent problem in which a hidden action has a persistent effect on output. We proposed a measure of deferral and, assuming the risk–averse agent has square root utility, characterized its comparative statics with respect to the degree of persistence and the common discount factor of the principal and the agent. We identified two countervailing effects that determine the optimal amount of deferral. The first effect is the intuitive one that by waiting to pay the agent until the principal has better information about the agent’s performance, incentives can be more efficiently implemented. The second effect is less intuitive, but also related to having better information. When the effect of effort is persistent over time this means that there is more than one output that provides information about hidden effort. This “better” information allows compensation to be less variable overall, and, sometimes, that lowers the expected payment in the second period made by the principal more than that of the first period. That is, in some cases, the second force dominates so deferred compensation can decrease as persistence increases.

Even in the absence of persistence, there is value to deferring compensation in multi–period incentive problems. To provide a benchmark, we studied our measure of deferral in the standard repeated moral hazard model. Interestingly, we found cases where there is more deferral of compensation in this model than in the corresponding persistence model. One significant difference between the two classes of models is that the comparative statics of deferral differed with respect to discounting.

¹³For a related discussion, see Jarque and Prescott (2015).

Persistent effects of inputs on output is a pervasive feature of production processes. While this class of models is applicable to a wide variety of problems, one particularly interesting application is the banking sector. Recent regulation attempts to limit bank risk by requiring banks to defer higher proportions of compensation than they have in the past. A first implication of our results is that while deferral is an important part of optimal contracting arrangements seeking to implement prudent risk-taking, more deferral is not always desirable and, somewhat surprisingly, more persistence in production does not necessarily mean that there should be more deferral. In that sense, mandatory minimums on deferral may very well be undesirable.

A second implication of our results concerns the measurement and assessment of incentives in compensation arrangements: assessment should not be based on the legal form of the compensation arrangements. That is, whether payments are labeled a bonus or a raise is pretty much irrelevant for incentives; what really matters is the actual state-contingent payments, regardless of their legal form, that are paid out. While this may seem like an obvious point, it is one that is often ignored. Our measure of deferral, which captures the time path of expected compensation, is independent of the legal form of payments and thus useful for assessing the timing of compensation.

References

- [1] Ábrahám, Á. and N. Pavoni. “Efficient Allocations with Moral Hazard and Hidden Borrowing and Lending: A Recursive Formulation,” *Review of Economic Dynamics*, Elsevier for the Society for Economic Dynamics, vol. 11(4), pages 781-803, October. (2008)
- [2] Agarwal, Sumit, and Itzhak Ben-David, “Loan Prospecting and the Loss of Soft Information.” (July 31, 2015). Charles A. Dice Center Working Paper No. 2012-7; Fisher College of Business Working Paper No. 2012-03-007. Available at SSRN: <http://ssrn.com/abstract=2046696>
- [3] Berg, Tobias and Puri, Manju and Rocholl, Jörg, Loan Officer Incentives, Internal Ratings, and Default Rates (September 2, 2014). AFA 2013 San Diego Meetings Paper. Available at SSRN: <http://ssrn.com/abstract=2022972>
- [4] Grossman, Sanford and Oliver D. Hart. “An Analysis of the Principal–Agent Problem.” *Econometrica* 51, Issue 1 (Jan.,1983), 7-46.
- [5] Cole, S., M. Kanz, L. Klapper. “Incentivizing Calculated Risk-Taking: Evidence from an Experiment with Commercial Bank Loan Officers.” *The Journal of Finance*, 70(2), 2015.

- [6] Federal Register, “Guidance on Sound Incentive Compensation Practices.” Vol. 75 (122), June 25, 2010.
- [7] Financial Crisis Inquiry Commission. “The financial crisis inquiry report: Final report of the National Commission on the Causes of the Financial and Economic Crisis in the United States.” Washington, DC. (2010)
- [8] Hartman-Glaser, B; Piskorski, T.; Tchisty, A. “Optimal Securitization with Moral Hazard.” *Journal of Financial Economics* (July 15, 2012).
- [9] Heider, F., and R. Inderst. “Loan Prospecting,” *Review of Financial Studies*, 25(8), 2381-2415 (2012)
- [10] Hertzberg, Andrew, Jose Maria Liberti, and Daniel Paravisini, “Information and Incentives Inside the Firm: Evidence from Loan Officer Rotation.” *Journal of Finance*, 65(3). 795-828.
- [11] Hoffmann, F., R. Inderst and M. Opp. “Regulating Deferred Incentive Pay”, mimeo (2013).
- [12] Holmström, B. “Moral Hazard and Observability,” *Bell Journal of Economics*, Vol. **10** (1) pp. 74-91. (1979)
- [13] Holmström, B. “Moral Hazard in Teams.” *Bell Journal of Economics* vol. 13 (Autumn), 1982, pp 324-340.
- [14] Hopenhayn, Hugo, and Arantxa Jarque. “Unobserved Persistent Productivity and Long Term Contracts.” *Review of Economic Dynamics*, vol. 13, 2010, pp 333-349.
- [15] Jarque, A., and E.S. Prescott. “Bank Compensation and Bank Risk Taking: The Organizational Economics View.” Richmond Fed WP 13-03 (2013)
- [16] N. Kocherlakota, Figuring Out the Impact of Hidden Savings on Optimal Unemployment Insurance, *Rev. Econ. Dynam.*, 7(3), 541-554 (2004).
- [17] Kim, S. K. “Efficiency of an Information System in an agency Model.” *Econometrica*, vol **63**(1), pages 89–102 (1995)
- [18] Kwon, I. “Incentives, Wages, and Promotions: Theory and Evidence.” *Rand Journal of Economics*, **37** (1), 100-120 (2006)
- [19] Mukoyama, T. and A. Sahin, “Repeated Moral Hazard with Persistence,” *Economic Theory*, vol. **25**(4), pages 831-854, 06 (2005)

- [20] Phelan, Christopher, “A Simple Model of Bank Employee Compensation.” Federal Reserve Bank of Minneapolis Working Paper 676, 2009.
- [21] Prescott, Edward Simpson and Robert M. Townsend. “Private Information and Intertemporal Job Assignments.” *Review of Economic Studies*, vol. 72, 2006, pp 531-548.
- [22] Rogerson, William P. “Repeated Moral Hazard.” *Econometrica*, Vol. 53, No. 1. (1985), pp 69-76.

9 Appendix 1: Derivations of optimal contract

9.1 Benchmark RMH

The FOCs are:

$$(u_i) : \frac{1}{U'(c_i)} = \lambda + \mu_1 (1 - LR_i), \quad i = l, h$$

$$(u_{ij}) : \frac{1}{U'(c_{ij})} = \lambda + \mu_1 (1 - LR_i) + \frac{\mu_{2i}}{\beta \text{Pr}_i} (1 - LR_j), \quad i, j = l, h$$

We can solve for the multipliers using the PC and the IC, and the fact that $1/U' = U/2$:

$$\lambda = \left(\frac{\bar{U}}{1 + \beta} + z(a) \right) \frac{1}{2}$$

$$\mu_1 = \frac{1}{2} \frac{z(a)}{v} \frac{1}{1 + \beta}$$

$$\mu_{2i} = \frac{1}{2} \beta \text{Pr}_i \frac{z(a)}{v},$$

where $v = \frac{(\pi_g - \pi_b)^2}{\pi_g(1 - \pi_g)}$ is the variance of the individual outcome realizations. Expressions for utilities:

$$u_i = \frac{1}{2} \left\{ \frac{\bar{U}}{1 + \beta} + z(a) + z(a) \frac{1 - LR_i}{v_1(1 + \beta)} \right\}$$

$$u_{ij} = \frac{1}{2} \left\{ \frac{\bar{U}}{1 + \beta} + z(a) + z(a) \left(\frac{1 - LR_i}{v(1 + \beta)} + \frac{1 - LR_j}{v} \right) \right\}$$

Given the expressions for the utilities, we can find expected consumption in each period:

$$E[c_1] = \lambda^2 + \mu^2 v_1$$

$$= \frac{1}{4} \left\{ \left(\frac{\bar{U}}{1 + \beta} + z(a) \right)^2 + \frac{[z(a)]^2}{v} \frac{1}{(1 + \beta)^2} \right\}$$

$$\begin{aligned}
E[c_2] &= \lambda^2 + \mu^2 v + \left(\frac{\mu_{2H}^2}{\beta^2 \pi} + \frac{\mu_{2L}^2}{\beta^2 (1 - \pi)} \right) v \\
&= \frac{1}{4} \left\{ \left(\frac{\bar{U}}{1 + \beta} + z(a) \right)^2 + \frac{[z(a)]^2}{v} \left(\frac{1}{(1 + \beta)^2} + 1 \right) \right\}.
\end{aligned}$$

9.2 Main model

For our utility specification, when the lower bound of utility is not binding¹⁴ we can solve for the closed form solution for the contingent utilities.¹⁵ From the FOCs,

$$(u_s) : \sqrt{c_s} = \lambda + \mu(1 - LR_s),$$

where λ is the multiplier of the IC constraint and μ is the multiplier of the IC constraint. Combining these with the PC and the IC, we have

$$\begin{aligned}
\lambda &= \frac{(\bar{U} + \phi)}{2(1 + \beta)} \\
\mu &= \frac{\phi}{2} \frac{1}{v_1 + \beta v_2},
\end{aligned}$$

so

$$u_s = 2\sqrt{c_s} = \frac{\bar{U} + \phi}{1 + \beta} + \phi \frac{1}{v_1 + \beta v_2} (1 - LR_s), \quad (24)$$

where $s \in \{l, h, ll, lh, hl, hh\}$. Note that $v_1 = v = \frac{(\pi_g - \pi_b)^2}{\pi_g(1 - \pi_g)}$.

Inverting $U(c) = 2\sqrt{c}$, we get that

$$\begin{aligned}
c_s &= [\lambda + \mu(1 - LR_s)]^2 \\
&= \frac{1}{4(1 + \beta)^2} \left[\bar{U} + \phi + \frac{\phi}{v} (1 - LR(y^t)) \right]^2.
\end{aligned}$$

¹⁴The necessary and sufficient condition on the parameters for the non-negativity constraint on square root utility not to be binding is:

$$\frac{\bar{U}}{a} + 1 \geq \max_m \frac{1}{\bar{v}^m} \left[\left(\frac{1 - \hat{\pi}^m}{1 - \pi^m} \right)^2 - 1 \right],$$

where π^m and $\hat{\pi}^m$ denote the probabilities of repayment on and off the equilibrium path, respectively, and \bar{v}^m the average discounted variance of the likelihood ratios, corresponding to model $m \in \{PB, PT, PV, SB, SV\}$.

¹⁵Note that if the domain constraint in (4) were to bind, some of the strict inequalities in Prop. 1 would become equalities.

Expected consumption at each t can be written as:

$$\begin{aligned}
E [c_t | a_h] &= E [[\lambda + \mu (1 - LR_s)]^2 | a_h] \\
&= \lambda^2 + \mu^2 v_t \\
&= \frac{1}{4(1 + \beta)^2} \left[(\bar{U} + \phi)^2 + \left(\frac{\phi}{v_1 + \beta v_2} \right)^2 v_t \right].
\end{aligned}$$

9.3 Extension

The constraints of the problem simplify to the PC

$$u_1 + \beta [\pi_g u_{2h} + (1 - \pi_g) u_{2l}] - \frac{\phi}{1 + \beta} - \beta \frac{\phi}{1 + \beta} = \bar{U}$$

and the second period ICs, which are both equal and hence simplify to a unique IC:

$$\begin{aligned}
\pi_g u_{2h} + (1 - \pi_g) u_{2l} - \frac{\phi}{1 + \beta} &= [\alpha \pi_g + (1 - \alpha) \pi_b] u_{2h} \\
&\quad + [1 - \alpha \pi_g - (1 - \alpha) \pi_b] u_{2l},
\end{aligned}$$

or, rearranging,

$$\begin{aligned}
\{\pi_g - [\alpha \pi_g + (1 - \alpha) \pi_b]\} u_{2h} + \{(1 - \pi_g) - [1 - \alpha \pi_g - (1 - \alpha) \pi_b]\} u_{2l} &= \frac{\phi}{1 + \beta} \\
(1 - \alpha) (\pi_g - \pi_b) (u_{2h} - u_{2l}) &= \frac{\phi}{1 + \beta}.
\end{aligned}$$

The first-order conditions of the problem are:

$$\begin{aligned}
(c_1) &: \frac{1}{U'(c_1)} = \sqrt{c_1} = \lambda \\
(c_j) &: \frac{1}{U'(c_i)} = \sqrt{c_i} = \lambda + \frac{\mu}{\beta} [1 - LR_j],
\end{aligned}$$

where λ is the multiplier of the IC constraint and μ is the multiplier of the IC constraint. For $U(c) = 2\sqrt{c}$, the PC together with the first-order conditions imply that

$$2\lambda + \beta \left[\pi_g 2 \left(\lambda + \frac{\mu}{\beta} [1 - LR_h] \right) + (1 - \pi_g) 2 \left(\lambda + \frac{\mu}{\beta} [1 - LR_l] \right) \right] - \phi = \bar{U}$$

$$\begin{aligned}
2\lambda + \beta 2\lambda + \beta 2 \frac{\mu}{\beta} \left[\pi_g - \frac{\pi_g \pi_b}{\pi_g} + 1 - \pi_g - \frac{(1 - \pi_g)(1 - \pi_b)}{(1 - \pi_g)} \right] - \phi &= \bar{U} \\
2(1 + \beta)\lambda - \phi &= \bar{U}
\end{aligned}$$

$$\lambda = \frac{\bar{U} + \phi}{2(1 + \beta)}.$$

The IC together with the first-order conditions imply that

$$\begin{aligned} (1 - \alpha) (\pi_g - \pi_b) 2 \left(\lambda + \frac{\mu}{\beta} [1 - LR_h] - \lambda - \frac{\mu}{\beta} [1 - LR_l] \right) &= \frac{\phi}{1 + \beta} \\ 2 \frac{\mu}{\beta} (1 - \alpha) (\pi_g - \pi_b) [LR_l - LR_h] &= \frac{\phi}{1 + \beta} \\ 2 \frac{\mu}{\beta} \frac{[(1 - \alpha) (\pi_g - \pi_b)]^2}{\pi_g (1 - \pi_g)} &= \frac{\phi}{1 + \beta} \\ \mu &= \frac{\phi \beta}{2(1 + \beta)} \frac{\pi_g (1 - \pi_g)}{[(1 - \alpha) (\pi_g - \pi_b)]^2}. \end{aligned}$$

Recall that $v = \frac{(\pi_g - \pi_b)^2}{\pi_g(1 - \pi_g)}$. In the extension with two efforts, the variance of the likelihood ratios in the second period, \tilde{v}_2 , will be a modified version of v :

$$\tilde{v}_2 = \frac{[(1 - \alpha) (\pi_g - \pi_b)]^2}{\pi_g (1 - \pi_g)} = (1 - \alpha)^2 v.$$

Hence, we can write

$$\mu = \frac{\phi}{2(1 + \beta)} \frac{\beta}{(1 - \alpha)^2} \frac{1}{v}.$$

$$\begin{aligned} u_1 &= \frac{\bar{U} + \phi}{(1 + \beta)} \\ u_{2i} &= \frac{\bar{U} + \phi}{(1 + \beta)} + \frac{\phi}{(1 + \beta)} \frac{\beta}{(1 - \alpha)^2} \frac{1}{v} (1 - LR_i), \text{ for } i = l, h. \end{aligned}$$

9.4 Implementation of optimal contract with real life schemes

9.4.1 Compensation package 1

In order to implement the optimal contract, this compensation scheme must satisfy the following system of equations:

$$\begin{aligned} c_l &= W + (1 - d_l) B_l \\ c_h &= W + (1 - d_h) B_h \\ c_{lh} &= W + d_l B_l \\ c_{hh} &= W + d_h B_h \\ c_l &= W + (1 - \theta_l) d_l B_l \\ c_h &= W + (1 - \theta_h) d_h B_h. \end{aligned}$$

Define the following notation for consumption spreads: $\Delta_0 = c_h - c_l$, Δ_h as the difference between the two highest second-period consumptions ($c_{hh} - c_{hl}$, when $\alpha > 0$), Δ_l as the difference between the two lowest second-period consumptions ($c_{lh} - c_{ll}$, when $\alpha > 0$), and $\Delta_{ll} = c_{ll} - c_l$. Using this notation, and noting that Proposition 1 implies that all these spreads are positive whenever $\alpha = 1$, this first scheme implies:

$$\begin{aligned}
W &= c_{ll} \\
B_h &= \Delta_{ll} + \Delta_0 + \Delta_h + \Delta_l \\
B_l &= \Delta_{ll} + \Delta_l \\
d_h &= \frac{\Delta_h + \Delta_l}{\Delta_{ll} + \Delta_0 + \Delta_h + \Delta_l} \\
d_l &= \frac{\Delta_l}{\Delta_{ll} + \Delta_l} \\
\theta_h &= \frac{\Delta_h}{\Delta_h + \Delta_l} \\
\theta_l &= 1.
\end{aligned}$$

That is, a bonus payment is given in the first period, with $B_h > B_l$, and a portion of each is deferred to the second period. Note that since $c_{hh} > c_h$ and $c_{lh} > c_l$, it follows that $d_l > 1/2$ and $d_h > 1/2$. Also, if output in the second period is low, all of the deferred part of bonus B_l is clawed back (that is, if the first period realization was high). On the other hand, only $\theta_h < 1/2$ is clawed back of the deferred part of B_h (that is, if the first-period realization was high).

When $\alpha = 0$, however, $c_{hh} = c_{lh}$ and $c_{hl} = c_{ll}$. Hence, there is no meaningful contingent spread in the second period, and the consumption spread is the same in both periods. Also, $\Delta_{ll} = 0$. This simple scheme is easily achievable with the following instruments:

$$\begin{aligned}
W &= c_l \\
B_h &= 2\Delta_0 \\
B_l &= 0 \\
d_h &= \frac{1}{2} \\
d_l &= 0 \\
\theta_h &= 0 \\
\theta_l &= 0.
\end{aligned}$$

That is, a bonus payment is given only after a high realization in the first period, of which 1/2 is deferred. Deferral is, then, lower when there is no persistence of effort. More importantly, because there is no information about effort contained in the second period, clawbacks in the second period are not useful. Hence, when using compensation packages of the type described

in scheme A, persistence calls for higher deferral than non-persistence, and persistence calls for clawbacks, while non-persistence does not.

9.4.2 Compensation package 2

In this alternative, the optimal contingent consumptions are implemented by the wage and bonus scheme as follows:

$$\begin{aligned}
 c_l &= W \\
 c_h &= W + B_0 \\
 c_{lh} &= W_l + B_l \\
 c_{hh} &= W_h + B_h \\
 c_{ll} &= W_l \\
 c_{hl} &= W_h.
 \end{aligned}$$

Using the notation for consumption spreads, this scheme simply implies, for $\alpha = 1$:

$$\begin{aligned}
 W &= c_l \\
 B_0 &= \Delta_0 \\
 W_l &= c_{ll} \\
 B_l &= \Delta_l \\
 W_h &= c_{hl} \\
 B_h &= \Delta_h,
 \end{aligned}$$

where, by Proposition 1, $W_l < W < W_h$ and $B_h < B_l$. When $\alpha = 0$, we simply have that $W_l = W_h = W$, and $B_l = B_h = B_0$. More importantly, under this compensation scheme there are no deferrals or clawbacks, *even when there is perfect persistence*.

10 Appendix 2: Proofs

Proof of Proposition 1. See Hopenhayn and Jarque (2010). ■

Proof of Corollary 1. The proof of Corollary 1 follows trivially from Proposition 1 and substitution the values of α in the corresponding likelihood ratios. ■

Proof of Lemma 1. Mukoyama and Sahin (2005) present a more general probability structure and then proceed to characterize the subset of probability and discount factor

combinations that imply perfect insurance in the first period. The conditions in Mukoyama and Sahin become, under our assumptions:

$$\begin{aligned}\beta \{ \pi_g - [\pi_g - \alpha (\pi_g - \pi_b)] \} &\geq \pi_g - [\alpha \pi_g + (1 - \alpha) \pi_b], \\ \beta \{ [\alpha \pi_g + (1 - \alpha) \pi_b] - \pi_b \} &\geq \pi_g - [\alpha \pi_g + (1 - \alpha) \pi_b].\end{aligned}$$

Simplifying, the first condition becomes:

$$\beta \geq \frac{1 - \alpha}{\alpha},$$

and the second condition becomes:

$$\beta \geq \frac{1 - \alpha}{\alpha}.$$

■

Proof of Proposition 2. See Rogerson (1985). ■

Proof of Lemma 2. We can write $v_1 = v$ as:

$$\begin{aligned}v_1 &= \sum_{i=l,h} f_1(r_j|a_h) \{LR_i - E[LR_i]\}^2 \\ &= \frac{(\pi_g - \pi_b)^2}{\pi_g(1 - \pi_g)} \\ &= \pi_b \frac{\pi_b}{\pi_g} + (1 - \pi_b) \frac{1 - \pi_b}{1 - \pi_a} - 1 \\ &= \hat{E}_1 - 1,\end{aligned}$$

where $\hat{E}_1 \equiv E_{t=1}[LR_i|a_l] > 1$ denotes the expectation of the likelihood ratios of the first period outcome realizations under the conditional distribution defined by a_l . The variance

in the second period can be written as

$$\begin{aligned}
v_2 &= \sum_{i=l,h, j=l,h} f_1(r_i|a_h)g_2(r_h|a_1) \{LR_{ij} - E[LR_{ij}]\}^2 \\
&= \sum_{ij} \Pr_i \Pr_j (LR_i LR_j - 1)^2 \\
&= \sum_{ij} \Pr_i \Pr_j (LR_i^2 LR_j^2 - 2LR_i LR_j + 1) \\
&= \sum_{ij} \Pr_i \Pr_j \left(\frac{\widehat{\Pr}_i^2 \widehat{\Pr}_j^2}{\Pr_i^2 \Pr_j^2} - 2 \frac{\widehat{\Pr}_i \widehat{\Pr}_j}{\Pr_i \Pr_j} + 1 \right) \\
&= \sum_{ij} \widehat{\Pr}_i \widehat{\Pr}_j \frac{\widehat{\Pr}_i \widehat{\Pr}_j}{\Pr_i \Pr_j} - 2 \sum_{ij} \widehat{\Pr}_i \widehat{\Pr}_j + \sum_{ij} \Pr_i \Pr_j \\
&= \sum_{ij} \widehat{\Pr}_i \widehat{\Pr}_j LR_i LR_j - 2 + 1 \\
&= \sum_i \widehat{\Pr}_i LR_i \left(\sum_j \widehat{\Pr}_j LR_j \right) - 1 \\
&= \sum_i \widehat{\Pr}_i LR_i \left[\hat{E}_2(\alpha) \right] - 1 \\
&= \hat{E}_2(\alpha) \sum_i \widehat{\Pr}_i LR_i - 1 \\
&= \hat{E}_1 \hat{E}_2(\alpha) - 1,
\end{aligned}$$

$$v_2 = \hat{E}_1 \hat{E}_2(\alpha) - 1,$$

where $\hat{E}_2(\alpha) \equiv E_{t=2}[LR_j|a_t]$ denotes the expectation of the likelihood ratios of the second period realization under the conditional distribution defined by a_t :

$$\hat{E}_2 = \pi_b \frac{\pi_b}{\alpha \pi_g + (1 - \alpha) \pi_b} + (1 - \pi_b) \frac{1 - \pi_b}{1 - [\alpha \pi_g + (1 - \alpha) \pi_b]}.$$

Clearly, given our assumption that $\pi_g > \pi_b$, we have that \hat{E}_2 is strictly greater than 1 if and only if $\alpha > 0$. Then, $\hat{E}_2 > 1$ implies $v_2 > v_1$. ■

Proof of Prop. 3. We are interested in showing that the effect of an increase in α on the deferral of pay could be positive or negative:

$$D = \frac{E[c_2]}{E[c_1]} = \frac{\lambda^2 + \mu^2 v_2}{\lambda^2 + \mu^2 v_1}$$

Recall that

$$v_2 = Var(LR_i LR_j) = \hat{E}_1 \hat{E}_2(\alpha) - 1,$$

where

$$\begin{aligned}\hat{E}_1 &= \pi_b \frac{\pi_b}{\pi_g} + (1 - \pi_b) \frac{(1 - \pi_b)}{(1 - \pi)} \\ \hat{E}_2(\alpha) &= \hat{p} \frac{\hat{p}}{p} + (1 - \hat{p}) \frac{(1 - \hat{p})}{(1 - p)}.\end{aligned}$$

Hence $\frac{\partial v_2}{\partial \alpha} = \hat{E}_1 \frac{\partial \hat{E}_2(\alpha)}{\partial \alpha}$ is positive if and only if $\frac{\partial \hat{E}_2(\alpha)}{\partial \alpha} > 0$. We can show that the conditional second-period variance, $v_2^c = \hat{E}_2(\alpha) - 1$, which is a simpler expression, increases with α . This will imply that \hat{E}_2 increases with α as well. We can write the conditional variance in the second period, $v_2^c \equiv Var[LR(r_2) | r_1, a_h]$, as

$$v_2^c = \frac{\alpha^2 (\pi_g - \pi_b)^2}{(\alpha \pi_g + (1 - \alpha) \pi_b) [1 - \alpha \pi_g + (1 - \alpha) \pi_b]},$$

or, using a notation shortcut $p \equiv \alpha \pi_g + (1 - \alpha) \pi_b$,

$$v_2^c = \frac{(p - \pi_b)^2}{p(1 - p)} = \hat{E}_2(\alpha) - 1.$$

Because there is no persistence of output, this conditional variance does not depend on the realization of r_1 .

$$\begin{aligned}v_2^c(\alpha) &= \hat{E}_2(\alpha) - 1 \\ &= \frac{(p - \hat{p})^2}{p(1 - p)} \\ &= \frac{[\alpha \pi + (1 - \alpha) \pi_b - \alpha \pi_b - (1 - \alpha) \pi_b]^2}{p(1 - p)} \\ &= \frac{\alpha^2 (\pi - \pi_b)^2}{p(1 - p)} \\ &= \frac{\alpha^2 (\pi - \pi_b)^2}{[\alpha \pi + (1 - \alpha) \pi_b] (1 - \alpha \pi - (1 - \alpha) \pi_b)} \\ &= \frac{\alpha^2 (\pi - \pi_b)^2}{[\alpha (\pi - \pi_b) + \pi_b] [1 - \pi_b - \alpha (\pi - \pi_b)]} \\ &= \frac{\alpha^2 (\pi - \pi_b)^2}{\alpha (\pi - \pi_b) (1 - \pi_b) - \alpha^2 (\pi - \pi_b)^2 + \pi_b (1 - \pi_b) - \alpha \pi_b (\pi - \pi_b)} \\ &= \frac{\alpha^2 (\pi - \pi_b)^2}{\pi_b (1 - \pi_b) - \alpha (\pi - \pi_b) (2\pi_b - 1) - \alpha^2 (\pi - \pi_b)^2}.\end{aligned}$$

Then

$$\begin{aligned}
\frac{\partial v_2^c(\alpha)}{\partial \alpha} &= \frac{(\pi - \pi_b)^2}{p^2(1-p)^2} \left\{ \begin{array}{l} 2\alpha [\pi_b(1 - \pi_b) - \alpha(\pi - \pi_b)(2\pi_b - 1) - \alpha^2(\pi - \pi_b)^2] \\ -\alpha^2 [-(\pi - \pi_b)(2\pi_b - 1) - 2\alpha(\pi - \pi_b)^2] \end{array} \right\} \\
&= \frac{(\pi - \pi_b)^2}{p^2(1-p)^2} \left\{ \begin{array}{l} 2\alpha [\pi_b(1 - \pi_b) - \alpha(\pi - \pi_b)(2\pi_b - 1) - \alpha^2(\pi - \pi_b)^2] \\ +\alpha^2 [(\pi - \pi_b)(2\pi_b - 1) + 2\alpha(\pi - \pi_b)^2] \end{array} \right\} \\
&= \frac{(\pi - \pi_b)^2}{p^2(1-p)^2} \alpha \left\{ \begin{array}{l} 2\pi_b(1 - \pi_b) - 2\alpha(\pi - \pi_b)(2\pi_b - 1) - 2\alpha^2(\pi - \pi_b)^2 \\ +\alpha(\pi - \pi_b)(2\pi_b - 1) + 2\alpha^2(\pi - \pi_b)^2 \end{array} \right\} \\
&= \frac{(\pi - \pi_b)^2}{p^2(1-p)^2} \alpha \{2\pi_b(1 - \pi_b) - \alpha(\pi - \pi_b)(2\pi_b - 1)\}.
\end{aligned}$$

The sign of the derivative will be determined by the sign of

$$2\pi_b(1 - \pi_b) - \alpha(\pi - \pi_b)(2\pi_b - 1).$$

Assume that $(\pi - \pi_b)(2\pi_b - 1) > 0$. Then the above expression takes its minimum value when $\alpha = 1$, and this value is positive:

$$\begin{aligned}
2\pi_b(1 - \pi_b) - (\pi - \pi_b)(2\pi_b - 1) &= 2\pi_b - 2\pi_b^2 - 2\pi\pi_b + \pi + 2\pi_b^2 - \pi_b \\
&= \pi_b - 2\pi\pi_b + \pi \\
&= \pi_b - \pi\pi_b + \pi - \pi\pi_b \\
&= \pi_b(1 - \pi) + \pi(1 - \pi_b) > 0.
\end{aligned}$$

If, instead, $(\pi - \pi_b)(2\pi_b - 1) < 0$, then the expression is obviously positive. So $\frac{\partial v_2^c(\alpha)}{\partial \alpha} > 0$. With this, we can go back to the expression for deferred pay:

$$D = \frac{\lambda^2 + \mu^2(\alpha)v_2(\alpha)}{\lambda^2 + \mu^2(\alpha)v_1}.$$

Since $\lambda = \frac{\bar{U} + a}{2(1+\beta)}$, we have that $\frac{\partial \lambda}{\partial \alpha} = 0$. Instead, since

$$\mu = \frac{\phi}{2} \frac{1}{v_1 + \beta v_2(\alpha)},$$

this multiplier changes with α . We can easily see that this derivative is negative:

$$\frac{\partial \mu}{\partial \alpha} = \frac{\phi}{2} \frac{-\beta \frac{\partial v_2}{\partial \alpha}}{[v_1 + \beta v_2(\alpha)]^2} < 0.$$

Combining the results above, we can see that $\frac{\lambda^2 + \mu^2 v_2}{\lambda^2 + \mu^2 v_1}$ may increase or decrease with α : on one hand, it may increase because $\lambda > 0$, $\mu > 0$ and $\frac{\partial v_2}{\partial \alpha} > 0$; on the other hand, it may

decrease because $\frac{\partial \lambda}{\partial \alpha} = 0$ and $\frac{\partial \mu}{\partial \alpha} < 0$ and $v_2 > v_1$. The formal derivative shows this:

$$\begin{aligned}
\frac{\partial D}{\partial \alpha} &= \frac{(2\mu \frac{\partial \mu}{\partial \alpha} v_2 + \mu^2 \frac{\partial v_2}{\partial \alpha}) (\lambda^2 + \mu^2 v_1) - (\lambda^2 + \mu^2 v_2) 2\mu \frac{\partial \mu}{\partial \alpha} v_1}{(\lambda^2 + \mu^2 v_1)^2} \\
&= \frac{\lambda^2 (2\mu \frac{\partial \mu}{\partial \alpha} v_2 + \mu^2 \frac{\partial v_2}{\partial \alpha} - 2\mu \frac{\partial \mu}{\partial \alpha} v_1) + \mu^2 [v_1 (2\mu \frac{\partial \mu}{\partial \alpha} v_2 + \mu^2 \frac{\partial v_2}{\partial \alpha}) - 2\mu v_2 \frac{\partial \mu}{\partial \alpha} v_1]}{(\lambda^2 + \mu^2 v_1)^2} \\
&= \mu \frac{\lambda^2 (2\frac{\partial \mu}{\partial \alpha} v_2 + \mu \frac{\partial v_2}{\partial \alpha} - 2\frac{\partial \mu}{\partial \alpha} v_1) + \mu [2\mu v_1 v_2 \frac{\partial \mu}{\partial \alpha} + \mu^2 \frac{\partial v_2}{\partial \alpha} v_1 - 2\mu v_1 v_2 \frac{\partial \mu}{\partial \alpha}]}{(\lambda^2 + \mu^2 v_1)^2} \\
&= \mu \frac{\lambda^2 (2\frac{\partial \mu}{\partial \alpha} (v_2 - v_1) + \mu \frac{\partial v_2}{\partial \alpha}) + \mu^3 \frac{\partial v_2}{\partial \alpha} v_1}{(\lambda^2 + \mu^2 v_1)^2}.
\end{aligned}$$

The sign of $\frac{\partial D}{\partial \alpha}$ will be determined by the sign of

$$\lambda^2 \left(2\frac{\partial \mu}{\partial \alpha} (v_2 - v_1) + \mu \frac{\partial v_2}{\partial \alpha} \right) + \mu^3 \frac{\partial v_2}{\partial \alpha} v_1.$$

This will be positive whenever

$$\left\| \lambda^2 2\frac{\partial \mu}{\partial \alpha} (v_2 - v_1) \right\| < \lambda^2 \mu \frac{\partial v_2}{\partial \alpha} + \mu^3 \frac{\partial v_2}{\partial \alpha} v_1,$$

or, since $v_2 > v_1$ but $\frac{\partial \mu}{\partial \alpha} < 0$,

$$2\lambda^2 \frac{\partial \mu}{\partial \alpha} (v_1 - v_2) < \mu \frac{\partial v_2}{\partial \alpha} (\lambda^2 + \mu^2 v_1).$$

■

Lemma 3 $\frac{\partial D}{\partial \alpha} > 0$ always at $\alpha = 0$.

Proof of Lemma 3. When $\alpha = 0$, we have that $v_2 = v_1 = v$, and hence $c_{lh} = c_{ll} = c_l$, as well as $c_{hl} = c_{hh} = c_{hh}$ (see Prop. 1). This implies $E[c_1] = E[c_2]$ and $D = 1$ when $\alpha = 0$. Consider $\alpha = \varepsilon$, where ε is a small but positive number. By Lemma 2, at $\alpha = \varepsilon$ we have $v_2 > v_1$, which implies $E[c_1] < E[c_2]$ and $D > 1$, and hence $\frac{\partial D}{\partial \alpha} > 0$. ■

Proof of Prop. 4. The derivative of the deferred pay for our main model with persistence is:

$$\begin{aligned}
\frac{\partial D}{\partial \beta} &= \frac{\frac{-2\phi^2 v_1^2 (1+v_1)^2}{(v_1+\beta v_2)^3 (1+\beta)^2} \left[(\bar{U} + \phi)^2 + \left(\frac{\phi}{v_1+\beta v_2} \right)^2 v_1 \right]}{\left[(\bar{U} + \phi)^2 + \left(\frac{\phi}{v_1+\beta v_2} \right)^2 v_1 \right]^2} \\
&\quad - \frac{\left[\left(\frac{\phi}{v_1+\beta v_2} \right)^2 v_1 (1+v_1) \right] \frac{-2\phi^2 v_1^2 (1+v_1)}{(v_1+\beta v_2)^3 (1+\beta)^2}}{\left[(\bar{U} + \phi)^2 + \left(\frac{\phi}{v_1+\beta v_2} \right)^2 v_1 \right]^2} \\
&= \frac{-2\phi^2 v_1^2 (1+v_1)^2 (\bar{U} + \phi)^2}{\left[(\bar{U} + \phi)^2 + \left(\frac{\phi}{v_1+\beta v_2} \right)^2 v_1 \right]^2 (v_1 + \beta v_2)^3 (1 + \beta)^2} < 0.
\end{aligned}$$

■

Proof of Prop. 5. Recall that $z(a_1) = z(a_2)$, and denote both by $z(a)$. The derivative of the deferred pay for the benchmark RMH is:

$$\frac{\partial D}{\partial \beta} = \frac{z(a)}{v} \frac{0 - \left\{ 2 \left\{ \frac{\bar{U}}{1+\beta} + z(a) \right\} \frac{-\bar{U}}{(1+\beta)^2} + \left(\frac{[z(a)]^2}{v} \right)^2 \left(\frac{-2}{(1+\beta)^3} \right) \right\}}{\left[\left(\frac{\bar{U}}{1+\beta} + z(a) \right)^2 + \frac{[z(a)]^2}{v} \frac{1}{(1+\beta)^2} \right]^2} > 0$$

Alternatively, we could use the expressions with the multipliers:

$$D = 1 + \frac{\left(\frac{\mu_{2h}^2}{\beta^2 \pi} + \frac{\mu_{2l}^2}{\beta^2 (1-\pi)} \right) v}{\lambda^2 + \mu_1^2 v}$$

$$\lambda = \frac{\frac{1}{2} \bar{U} + z(a_1) + \beta z(a_2)}{1 + \beta} = \frac{\bar{U} + z(a) (1 + \beta)}{2 (1 + \beta)};$$

$$\frac{\partial \lambda}{\partial \beta} = \frac{-\bar{U}}{2 (1 + \beta)^2} < 0$$

$$\begin{aligned}
\mu_1 &= \frac{1}{2} \frac{z(a_1)}{v} \frac{1}{1 + \beta}; \\
\frac{\partial \mu_1}{\partial \beta} &= \frac{1}{2} \frac{z(a_1)}{v} \frac{-1}{(1 + \beta)^2} < 0
\end{aligned}$$

$$\begin{aligned}\mu_{2i} &= \frac{1}{2}\beta \Pr_i \frac{z(a_2)}{v}; \\ \frac{\partial \mu_{2i}}{\partial \beta} &= \frac{1}{2} \Pr_i \frac{z(a_2)}{v} > 0\end{aligned}$$

However, $\frac{\mu_{2h}^2}{\beta^2 \pi} = \frac{\mu_{2l}^2}{\beta^2 (1-\pi)} = \frac{[z(a_2)]^2}{2^2 v}$, and the derivative of that term with respect to β is zero, so only the denominator of D changes (decreases) when β increases. In other words, changes in the discount factor do not affect expected consumption in the second period, but they decrease expected consumption in the first. Hence, we confirm that $\frac{\partial D^{RMH}}{\partial \beta} > 0$. ■

Proof of Claim in Example 3. Note that, for X and Y independent random variables, it can be shown (Goodman, 1962) that:

$$Var(XY) = (E[Y])^2 Var(X) + (E[X])^2 Var(Y) + Var(X) Var(Y).$$

In our main model with persistence, the second-period likelihood ratio, LR_{ij} , is equal to the product of the likelihood ratios of the two individual outputs, and $E[LR_i|a_h] = E[LR_j|a_h] = 1$. Hence,

$$\begin{aligned}v_2 &= Var(LR_{ij}) \\ &= Var(LR_i LR_j) \\ &= 1^2 v_1 + 1^2 v_2 + v_1 v_2.\end{aligned}$$

Recall that $v_1 = v$. Note also that when $\alpha = 1$, we have that $v_2 = v_1(1 + v_1)$, so $v_2 = v(2 + v)$. Substituting our assumptions of $\bar{U} = 0$ and $z(a_1) = z(a_2) = \frac{\phi}{1+\beta}$ in equations 21 and 23, we have that $D^{RMH} > D^{PB}$ if and only if

$$\begin{aligned}1 + \frac{\frac{\phi^2}{(1+\beta)^2 v}}{\frac{\phi^2}{(1+\beta)^2} + \frac{\phi^2}{(1+\beta)^2 v} \frac{1}{(1+\beta)^2}} &> \frac{\phi^2 + \frac{\phi^2}{[v+\beta v(2+v)]^2} v(2+v)}{\phi^2 + \frac{\phi^2}{[v+\beta v(2+v)]^2} v} \\ \frac{1 + \frac{1}{v(1+\beta)^2} + \frac{1}{v}}{1 + \frac{1}{v(1+\beta)^2}} &> \frac{1 + \frac{1}{[v+\beta v(2+v)]^2} v(2+v)}{1 + \frac{1}{[v+\beta v(2+v)]^2} v} \\ \frac{\frac{v(1+\beta)^2 + 1 + (1+\beta)^2}{v(1+\beta)^2}}{\frac{v(1+\beta)^2 + 1}{v(1+\beta)^2}} &> \frac{\frac{[v+\beta v(2+v)]^2 + v(2+v)}{[v+\beta v(2+v)]^2}}{\frac{[v+\beta v(2+v)]^2 + v}{[v+\beta v(2+v)]^2}} \\ \frac{v(1+\beta)^2 + 1 + (1+\beta)^2}{v(1+\beta)^2 + 1} &> \frac{[v + \beta v(2+v)]^2 + v(2+v)}{[v + \beta v(2+v)]^2 + v}.\end{aligned}$$

Substituting our assumption of $\beta = 1$,

$$\begin{aligned}
\frac{4v+1+4}{4v+1} &> \frac{[v+v(2+v)]^2+v(2+v)}{[v+v(2+v)]^2+v} \\
\frac{4v+5}{4v+1} &> \frac{[3v+v^2]^2+2v+v^2}{[3v+v^2]^2+v} \\
\frac{4v+5}{4v+1} &> \frac{9v^2+6v^3+v^4+2v+v^2}{9v^2+6v^3+v^4+v} \\
\frac{4v+5}{4v+1} &> \frac{v^3+6v^2+10v+2}{v^3+6v^2+9v+1} \\
(4v+5)(9v+6v^2+v^3+1) &> (4v+1)(10v+6v^2+v^3+2) \\
36v^2+24v^3+4v^4+4v+45v+30v^2+5v^3+5 &> 40v^2+24v^3+4v^4+8v+10v+6v^2+v^3+2 \\
4v^4+29v^3+66v^2+49v+5 &> 4v^4+25v^3+46v^2+18v+2 \\
29v^3+66v^2+49v+5 &> 25v^3+46v^2+18v+2 \\
4v^3+20v^2+31v+3 &> 0.
\end{aligned}$$

■

Proof of Prop. 6. In the extension to two actions, we can write the variance in the second period as:

$$v_2^B = (1 - \alpha)^2 v.$$

Hence, this variance will decrease when α increases:

$$\frac{\partial v_2^B}{\partial \alpha} = -2(1 - \alpha)v < 0.$$

We also can solve for the multiplier of the IC in (17) as a function of v :

$$\mu = \beta \frac{\phi}{v_2^B} = \beta \frac{\phi}{(1 - \alpha)^2 v}.$$

We will have that μ increases with α .

$$\frac{\partial \mu}{\partial \alpha} = \frac{\phi}{v(1 - \alpha)^3} > 0.$$

Whether deferral D increases or decreases with α will depend on which of these two forces dominates:

$$D = \frac{\lambda^2 + \mu^2 v_2^B}{\lambda^2}.$$

We can substitute for μ and v_2^B as a function of α :

$$\begin{aligned} D &= \frac{\lambda^2 + \left(\frac{\beta\phi}{(1-\alpha)^2 v}\right)^2 (1-\alpha)^2 v}{\lambda^2} \\ &= \frac{\lambda^2 + \frac{\beta^2 \phi^2}{(1-\alpha)^2 v}}{\lambda^2} \end{aligned}$$

Taking the partial with respect to α , we find that it is always positive:

$$\frac{\partial D}{\partial \alpha} = \frac{2\beta^2 \phi^2}{\lambda^2 (1-\alpha)^3 v} > 0.$$

■