

Equilibrium Contracts and Boundedly Rational Expectations*

Heiner Schumacher[†] Heidi Christina Thysen[‡]

Version: October 14, 2020

Abstract

We study a principal-agent framework in which the agent forms beliefs based on a misspecified subjective model of the principal's project. She fits this model to the objective probability distribution to predict output under alternative actions. Misspecifications in the subjective model may lead to biased beliefs. However, under mild restrictions, the agent has correct beliefs on the equilibrium path so that the optimal contract is non-exploitative. This allows for a behavioral version of the informativeness principle: The optimal contract conditions on an additional variable only if it is informative about the action according to the agent's subjective model. We further characterize when misspecifications affect the optimal contract. One implication of this characterization is that the scope for belief biases depends on the agent's job, e.g., her position in the hierarchy.

Keywords: Bayesian Networks, Principal-Agent Relationship, Bounded Rationality

JEL Classification: D03, D82, D86

*We gratefully acknowledge financial support by the ERC Advanced Investigator grant no. 692995. We thank Yair Antler, Felix Bierbrauer, Markus Dertwinkel-Kalt, Kfir Eliaz, Florian Englmaier, Guido Friebel, Paul Heidhues, Johannes Johnen, Heiko Karle, Michael Kosfeld, Botond Kőszegi, Gilat Levy, Debraj Ray, Ronny Razin, Karl Schlag, Klaus Schmidt, Dirk Sliwka, Balázs Szentés, and Yves Le Yaouanq, as well as seminar audiences at Aarhus University, Boston University, University of Cambridge, University of Cologne, Dalhousie University, Humboldt-University Berlin, University of Konstanz, London School of Economics, Ludwig-Maximilians University Munich, Tilburg University, the EEA-ESEM 2017 in Lisbon, the ESSET 2018 in Gerzensee, and the Psychology and Economics of Causal Reasoning Workshop 2019 in London for their valuable comments and suggestions. We especially thank Ran Spiegler for his invaluable support of the project. The usual disclaimer applies.

[†]KU Leuven and University of Innsbruck. E-mail: heiner.schumacher@kuleuven.be.

[‡]London School of Economics. E-mail: h.c.thysen@lse.ac.uk.

1 Introduction

The canonical principal-agent model of contracting under asymmetric information assumes that the agent knows the probabilistic consequences of all available actions. Formally, these are defined by a production function $p(y | a)$, where y is the contractible output and a the agent's action. Given the incentives provided by the contract, the agent chooses an action that – according to this function – maximizes her expected payoff. However, in an organization, $p(y | a)$ is typically a complex object. It may reflect information that is unavailable to the agent or that the agent cannot process due to cognitive limitations. Herbert Simon therefore proposed that administrative behavior must be “boundedly rational” (Simon 1947, 1955).

There are now many well-documented empirical cases where agents choose inferior actions, despite substantial experience and incentives for behavioral change. Bloom et al. (2013) test the impact of basic management practices, such as inventory control or the measurement of quality defects, on productivity in a number of large Indian textile factories. The factory managers often believed that these practices would not improve profits. Yet, their introduction increased productivity by 17 percent on average, and many practices were adopted later on after substantial consulting. This and many other cases¹ show that even experienced decision makers may ignore important aspects of their operation.

In this paper, we examine a contracting framework in which the agent exhibits misspecifications in her thinking. The classic approach to contracting with boundedly rational agents is to assume directly that beliefs $\hat{p}(y | a)$ about the production function are biased so that $\hat{p}(y | a) \neq p(y | a)$. An important implication of this approach is that the optimal contract may exploit the agent, in the sense that her (true) expected payoff falls below her reservation utility (e.g., Kőszegi 2014). However, it is unclear how sustainable biased beliefs – and hence exploitation – would be when the agent gathers experience. In the example above, the factory managers most likely know the expected outcomes from their usual actions, they just do not anticipate correctly what would happen if they change their behavior.

We therefore apply a new approach and assume that the agent estimates $p(y | a)$ based on data generated by the true production process, the implemented strategy α^* , and a non-parametric subjective model \mathcal{R} . A model \mathcal{R} is a collection of variables and causal relationships between these variables. It captures what the agent knows about the production process. This model may be misspecified. For example, it may be “too simple” relative to the complexity of

¹Nuland (2004) describes that, until the invention of germ theory, doctors saw no value in washing hands before treating patients, thereby killing many individuals through the transmission of diseases. This did not change even after being confronted with clear statistical evidence for the effectiveness of hygienic measures. Hanna et al. (2014) find that experienced farmers ignore important input dimensions and thus produce off the Pareto frontier. Blake et al. (2015) analyze a case where tech companies greatly overestimate the effectiveness of their online marketing efforts by neglecting the relationship between search clicks and purchase intent.

the organization: Empirical regularities that matter for the principal’s project may not appear in \mathcal{R} . The agent’s subjective beliefs about $p(y | a)$ will be denoted by $p_{\mathcal{R}}(y | a; \alpha^*)$. An equilibrium contract implements a strategy α^* if it is optimal for the agent to follow α^* under this contract given her beliefs $p_{\mathcal{R}}(y | a; \alpha^*)$. We study the properties of the optimal equilibrium contract, and obtain several new results on optimal contracting and organization that we would not get (or get only under very specific assumptions) if we directly choose beliefs $\hat{p}(y | a)$.

To capture the agent’s limited understanding of her environment, we use Spiegler’s (2016) Bayesian network approach. As an illustration, consider the following example:

“Marketer Example.” The agent is a marketer whose job is to increase sales y . One strategy to increase sales is to make cold-calls $a \in \{0, 1\}$, that is, calling potential customers without prior consent. Making cold-calls improves consumers’ information $x_1 \in \{0, 1\}$ about the firm’s product, but also reduces the firm’s reputation $x_2 \in \{0, 1\}$ since some customers are annoyed by being cold-called. Expected sales increase both in consumer information x_1 and reputation x_2 . However, when choosing her action, the marketer does not take the firm’s reputation into account. The only mechanism on her mind is that making cold-calls improves consumer information, and that more information translates into more sales.²

The Bayesian network approach roughly works as follows in the marketer example.³ The setting describes an “extended production function” $p(x_1, x_2, y | a)$, i.e., a joint probability distribution over the realization of consumer information, reputation and sales for any given action. This function captures the objective model \mathcal{R}^* of the project: \mathcal{R}^* contains all relevant variables, {action, consumer information, reputation, sales}, and the causal relationships between these variables. The agent’s subjective model \mathcal{R} is a simplified version of \mathcal{R}^* as it only contains the variables {action, consumer information, sales}, and their causal relationships. Her beliefs are derived by fitting \mathcal{R} to the objective probability distribution, which is generated by the implemented strategy α^* and the extended production function $p(x_1, x_2, y | a)$. Thus, the different elements in the agent’s subjective model \mathcal{R} are quantified using input from the true data-generating process. Combining these elements yields the agent’s subjective beliefs $p_{\mathcal{R}}(y | a; \alpha^*)$, which in general are not invariant to changes in α^* .

²A related mechanism is “demarketing” (Miklós-Thal and Zhang 2013): Extensive marketing can backfire since it may be a signal for low quality. One recent example are credit cards from non-financial companies that offer exceptionally large cash backs, but allow the issuing company to access (and exploit) consumers’ transaction histories (Kominers 2017). The latter aspect is usually hidden in the fine-print, and even sales representatives are often unaware of it. Nevertheless, some savvy consumers who are concerned with privacy may reject offers that just seem to be too generous.

³Missing technical details will be explained thoroughly in the next section.

We show that the optimal equilibrium contract exhibits the following features. First, a weak restriction on the agent’s subjective model guarantees that the participation constraint is not affected. This restriction is that \mathcal{R} is “perfect”, which means that the agent takes into account the link between any two variables in \mathcal{R} that have a joint influence on a third variable in \mathcal{R} . She then correctly predicts the marginal equilibrium distribution over output (Spiegler 2017), so that the optimal equilibrium contract does not exploit the agent. Importantly, a perfect \mathcal{R} ensures in many cases that there are no informational cues in the data the agent gathers on the equilibrium path that could alert her about the misspecification in \mathcal{R} .

Second, the principal may strictly benefit from the misspecification in the agent’s model even when exploitation is infeasible. In the marketer example, if the principal implements making cold-calls, then, by not taking reputation into account, the agent overestimates the drop in sales after deviation to not making cold-calls, i.e., she is “control optimistic” as defined by Spinnewijn (2013). This relaxes the incentive compatibility constraint, so that the principal can implement cold-calls with fewer incentives than if the agent had rational expectations.

Third, when \mathcal{R} is perfect, the incentive scheme in the optimal equilibrium contract appears to the agent as optimal for the principal. The agent then cannot deduct from the shape of incentives that her beliefs are biased. This is again different from the optimal contract under exogenously fixed biased beliefs where the agent may notice that the principal is betting against her. We show that in many cases the optimal equilibrium contract is “justifiable”, i.e., it is optimal for the principal from the agent’s point of view.

Taken together, these results show that an agent’s misperceptions can be sustainable in an organizational context: Neither her experiences on the equilibrium path nor the shape of the incentive contract inform the agent about the mistake in her thinking, and the principal benefits from this mistake. Building on these insights, we further analyze three topics in contract theory and organizational economics that are difficult to address with exogenously given beliefs: First, we derive a behavioral version of the informativeness principle. Second, we characterize when misspecifications in the agent’s model affect her beliefs. And third, we revisit some classic comparative statics results. We briefly describe each topic in turn.

An important question in contract theory is on which variables the optimal contract should condition the agent’s wage. According to the informativeness principle (e.g., Holmström 1979, Chaigneau et al. 2019), the optimal contract conditions on an additional signal z only if z provides information about the agent’s action that is not contained in y . We can derive an analogous statement when the agent has correct expectations on the equilibrium path about the joint distribution of y and z (with a further qualification this holds if \mathcal{R} is perfect). In this case, the optimal equilibrium contract conditions on z only if the agent’s action a and z are not independent conditional on y according to the agent’s subjective beliefs. This result does not

depend on other properties of the agent’s subjective model \mathcal{R} , and hence would hold in any setting where the agent’s beliefs about the joint distribution of y and z are correct. Nevertheless, we can use results from the Bayesian network literature to state sufficient conditions on \mathcal{R} so that the result’s requirements are satisfied. We apply these findings to provide a new explanation for why most executive compensation contracts do not condition on peer-performance (Bertrand and Mullainathan 2001, Bebchuk and Fried 2004).

Next, misspecifications in \mathcal{R} do not always affect the agent’s beliefs and optimal equilibrium contract. The agent is “behaviorally rational” if she correctly anticipates the production function, or, formally, $p_{\mathcal{R}}(y | a; \alpha) = p(y | a)$ for all possible a and α , regardless of the parametrization of the extended production function. We can find a correspondence $H^*(\mathcal{R}^*)$ which indicates for a given objective model \mathcal{R}^* the set of variables the agent must take into account in her simplified subjective model \mathcal{R} so that she is behaviorally rational. We show that $H^*(\mathcal{R}^*)$ is often a strict subset of the variables in \mathcal{R}^* , and that the difference between a variable $i \in H^*(\mathcal{R}^*)$ and a variable $j \notin H^*(\mathcal{R}^*)$ can be quite nuanced.

The characterization of $H^*(\mathcal{R}^*)$ shows which variables matter for the agent’s beliefs. An important interpretation of the objective model \mathcal{R}^* is that it captures the agent’s job, i.e., through which tasks, interactions, and decision-making powers she influences the final output. We can have two extended production functions that give rise to the same “reduced-form” production function $p(y | a)$, but that differ in their causal model \mathcal{R}^* , and hence in the extent to which simplifications affect $p_{\mathcal{R}}(y | a; \alpha^*)$. This allows us to examine which organizational features potentially cause the agent to overestimate the productivity of her effort. Consider an agent in a management position in which her effort influences the behavior of other workers (e.g., a group of marketers). If the agent does not understand the difficulties of their job (e.g., that cold-calls have a partial negative effect on sales through their effect on firm reputation), she overestimates her subordinates’ – and hence her own – productivity. There are different instances where this could happen: The agent may be a technical expert who is promoted into a management position in which she oversees the actions of workers whose job she does not fully understand. Alternatively, it may be the case that subordinates do not communicate the problems they face to their managers (due to career concerns). These phenomena are usually discussed critically in the management literature, but in our framework they advance the agent’s effort motivation and hence benefit the principal.

Finally, our framework allows for comparative statics since the agent’s beliefs are derived from the parameters of the true production process. We consider two classic comparative static results from the contract theory literature: the trade-off between risk and incentives, and the relationship between team size and incentives. For both cases, we show that the original results may no longer hold when the agent’s subjective model is misspecified.

Related Literature. Our basic model is the principal-agent framework introduced by Holmström (1979) and Grossman and Hart (1983). Holmström (1979) states a version of the informativeness principle. A generalization of it can be found in, e.g., Chaigneau et al. (2019). In the canonical framework, both principal and agent know the production function $p(y | a)$.

There are different approaches in behavioral contract theory that relax the assumption of unbiased beliefs about $p(y | a)$. First, several contracting models directly assume that the agent's beliefs about the production function are biased, i.e., $\hat{p}(y | a) \neq p(y | a)$; see Fang and Moscarini (2005), Van den Steen (2005), Gervais and Goldstein (2007), Santos-Pinto (2008), De la Rosa (2011), Sautmann (2007, 2013), Spinnewijn (2013, 2015). Specifically, this approach is used to model an overconfident agent who overestimates the probability of good states and underestimates the probability of bad states. This typically allows the principal to exploit the agent by paying more after high output and much less after low output, in which case the agent's expected payoff is below her reservation utility.

Second, a rich literature builds state-space models of “unawareness” (e.g., Dekel et al. 1998, Heifetz et al. 2006, 2013) and applies them to contracting settings. Auster (2013) examines a principal-agent model with an agent who is unaware of some output levels y , which again implies that the contract is exploitative. Von Thadden and Zhao (2012, 2014) assume that the agent is unaware of her available actions a and chooses a default action unless the principal educates her; unawareness then relaxes incentive compatibility at the default action.

Third, in order to justify biased beliefs, several papers assume that the agent knows the link between action and outcomes $p(y | a)$, but potentially gains from holding biased beliefs. She then chooses beliefs $\hat{p}(y | a)$ that solve the trade-off between the losses from biased decision-making and the gains from managing a self-control problem (Bénabou and Tirole 2002) or from enjoying anticipatory utility (Brunnermeier and Parker 2005, Kőszegi 2006). For an organizational context, Bénabou (2013) shows how the interaction between group members can make the suppression of bad news a strategic complement, so that collective denial of adverse signals (“groupthink”) occurs in equilibrium. Immordino et al. (2015) show that if the anticipatory utility is not too important, the principal may provide incentives so that it is optimal for the agent to choose correct beliefs.

Our approach to boundedly rational expectations and contracting is more conservative. The agent derives her beliefs from the true data-generating process, as in the canonical model; she just may not take into account all empirical regularities that matter for the principal's project. The misspecification in the agent's subjective model may cause her to overestimate her productivity, but, under a weak restriction, she still correctly anticipates the equilibrium distribution over output.

We also contribute to the literature on Bayesian networks/directed acyclic graphs (DAGs),

which have been used extensively in the artificial intelligence literature. Pearl (2009) promotes the view that DAGs represent causal relationships and provides a broad introduction to DAGs. In economics, Spiegel (2016, 2017) uses Bayesian networks to model agents with boundedly rational expectations. DAGs provide a general method to capture a variety of different inference errors such as reverse causation and coarseness. We build on these insights and apply them to contracting. Other recent papers use causal models to capture boundedly rational decision makers in monetary policy (Spiegel 2019), political competition (Eliaz and Spiegel 2020), Bayesian persuasion (Eliaz et al. 2019), and decision theory (Schenone 2020).

The remainder of the paper is organized as follows. Section 2 describes our framework. In Section 3, we examine how a misspecification in the agent’s subjective model affects the optimal contract. In Section 4, we state a behavioral version of the informativeness principle. In Section 5, we characterize when a misspecification leads to biased beliefs about the production function, and illustrate the implications of this characterization. In Section 6, we revisit two classic comparative statics of the canonical contracting framework. Section 7 concludes. Proofs and further results can be found in the appendix and a supplementary appendix.

2 The Model

We consider a standard principal-agent problem and combine it with the Bayesian network model of boundedly rational beliefs, as introduced in Spiegel (2016).

Basic Framework. Let $A \subset \mathbb{R}$ be a finite set of actions, $Y \subset \mathbb{R}$ a finite set of outputs, and $W \subseteq \mathbb{R}^{|Y|}$ the set of possible incentive schemes. The principal proposes a contract $(w(y), q(a))$, where $w(y) \in W$ is the agent’s wage conditional on the output $y \in Y$ and $q(a) \in \Delta(A)$ is the probability with which the principal wishes the agent to choose action $a \in A$. The agent can reject or accept the contract. If she rejects it, she enjoys the outside option value \bar{U} , while the principal earns zero. If she accepts the contract, she chooses an action $a \in A$. The agent’s personal cost of choosing a is given by a function $c(a)$. The action stochastically influences the project’s output. The agent’s utility from wage w is given by the utility function $u(w)$, which is weakly concave and exhibits $\lim_{w \rightarrow -\infty} u(w) = -\infty$. When the output is y and the agent’s action is a , the principal’s payoff is $V = y - w(y)$ and the agent’s payoff is $U = u(w(y)) - c(a)$.

Causal Structure. We model the causal structure through which the agent’s action affects the output. Let $N^* = \{0, \dots, n\}$ be the set of relevant variables (or nodes). This set contains the agent’s action and output, but may also include other variables. A generic realization of variable i is given by $x_i \in X_i$, where X_i is a finite set that contains at least two elements. Node

0 is the agent's action ($x_0 = a$, $X_0 = A$) and node n is the output ($x_n = y$, $X_n = Y$). The state is a vector $x^* = (x_0, x_1, \dots, x_n)$ and the set of all states is $X^* = \times_{i \in N^*} X_i$. For every subset $M \subseteq N^*$ and $x^* \in X^*$, let $x_M = (x_k)_{k \in M}$.

Denote by $p(x_1, \dots, x_n \mid a)$ the extended production function. For any action $a \in A$, it has full support over $X_1 \times \dots \times X_n$. We represent its causal structure by an irreflexive, asymmetric, and acyclic binary relation R^* over N^* , and denote it by the DAG $\mathcal{R}^* = (N^*, R^*)$, see the graph on the left of Figure 1 for an example. For two nodes $i, j \in N^*$ one may read iR^*j as “node i impacts on node j .” The set of nodes that influence i is defined, with abuse of notation, as $R^*(i) = \{j \in N^* \mid jR^*i\}$. Nothing influences the agent's action, $R^*(0) = \emptyset$. The probability distribution over states, $p(x^*) \in \Delta(X^*)$, then naturally factorizes according to \mathcal{R}^* via the formula

$$p(x^*) = q(x_0) \prod_{i \in N^* \setminus \{0\}} p(x_i \mid x_{R^*(i)}). \quad (1)$$

The “objective model” \mathcal{R}^* is one of the sparsest DAGs so that $p(x^*)$ factorizes according to \mathcal{R}^* . That is, \mathcal{R}^* contains exactly those conditional independence assumptions that are satisfied by $p(x^*)$.⁴



Figure 1: An objective model \mathcal{R}^* (left) and the agent's subjective model \mathcal{R} (right).

Beliefs, Personal Equilibrium, and Equilibrium Contract. The agent has her own subjective model $\mathcal{R} = (N, R)$, see the graph on the right of Figure 1 for an example. We assume that $\{0, n\} \in N \subseteq N^*$ and $R(0) = \emptyset$. The assumption that the agent includes her own action and the output in her subjective model ensures that her utility is measurable with respect to her beliefs. $N \subseteq N^*$ is assumed purely for simplicity. $R(0) = \emptyset$ implies that the agent knows that she does not receive any information about other variables prior to choosing an action, and that she has correct beliefs about the marginal distribution over her own action.

Definition 1. We say that \mathcal{R} is misspecified if $\mathcal{R} \neq \mathcal{R}^*$, and that \mathcal{R} is a simplification if $N \subset N^*$ and $R = N \times N \cap \mathcal{R}^*$.

⁴This rules out trivial cases such as when the objective distribution is consistent with the agent's DAG (as defined below), but the agent's DAG excludes links that are in \mathcal{R}^* .

A simplification is a misspecification where the agent's subjective model \mathcal{R} emerges from \mathcal{R}^* by dropping nodes from \mathcal{R}^* and the links adjacent to them. It will receive considerable attention in this paper, but most of our main results do not depend on whether the misspecification is a simplification or not. Denote by $x = (x_i)_{i \in N}$ the state vector for the agent's subjective model and $X = \times_{i \in N} X_i$. The agent fits her subjective model \mathcal{R} to the data generated by $p(x^*)$, so her beliefs factorize according to the formula

$$p_{\mathcal{R}}(x) = q(x_0) \prod_{i \in N \setminus \{0\}} p(x_i | x_{R(i)}). \quad (2)$$

Thus, all the conditional independence assumptions embedded in \mathcal{R} also appear in the agent's beliefs. For example, when the agent's subjective model is \mathcal{R} from Figure 1, her beliefs factorize according to $p_{\mathcal{R}}(a, x_1, y) = q(a)p(x_1 | a)p(y | x_1)$, where $q(a)$, $p(x_1 | a)$ and $p(y | x_1)$ follow from the probability distribution $p(x^*)$. Given the objective model in Figure 1, $p(y | x_1)$ will depend on $q(a)$ through variables 2 and 3. Hence, the agent's beliefs in general depend on $q(a)$. We therefore augment notation to indicate which strategy $q(a)$ is used when deriving beliefs and write $p_{\mathcal{R}}(x; q(a))$ instead of $p_{\mathcal{R}}(x)$. For any subset $M \subset N$, the agent's belief about the marginal distribution over x_M is $p_{\mathcal{R}}(x_M; q(a)) = \sum_{x_{N \setminus M} \in X_{N \setminus M}} p_{\mathcal{R}}(x_M, x_{N \setminus M}; q(a))$.

The agent follows the prescribed strategy from the contract only if it maximizes her expected utility given the wage scheme $w(y)$ and her subjective beliefs about the output conditional on her action, which we denote by $p_{\mathcal{R}}(y | a; q(a))$. These are computed as

$$p_{\mathcal{R}}(y | a; q(a)) = \frac{p_{\mathcal{R}}(a, y; q(a))}{\sum_{y \in Y} p_{\mathcal{R}}(a, y; q(a))}. \quad (3)$$

To close the model, we need to specify the agent's strategy $q(a)$ that is used to derive these beliefs. We adapt the personal equilibrium concept from Spiegler (2016) to our setting.

Definition 2. *The strategy $q(a)$ is a personal equilibrium at \mathcal{R} and $w(y)$ if for all actions $a \in A$ in the support of $q(a)$ we have*

$$a \in \arg \max_{a'} \sum_{y \in Y} p_{\mathcal{R}}(y | a'; q(a)) u(w(y)) - c(a'),$$

where $p_{\mathcal{R}}(y | a'; q(a)) = \lim_{k \rightarrow \infty} p_{\mathcal{R}}(y | a'; p^k(a))$ for all actions $a' \in A$ and a sequence $p^k(a) \rightarrow q(a)$ of fully mixed strategy profiles.

With the full support assumption, a fully mixed action profile ensures that all conditional probabilities are well-defined. The definition requires that equilibrium beliefs are the limit of a sequence of fully mixed profiles. A personal equilibrium always exists in our framework; see

Appendix A.1. We call a contract $(w(y), q(a))$ an “equilibrium contract” if $q(a)$ is a personal equilibrium at \mathcal{R} and $w(y)$. An optimal equilibrium contract is an equilibrium contract that maximizes the principal’s expected payoff. For convenience, we denote beliefs by $p_{\mathcal{R}}(y | a; a^*)$ when a pure action a^* is implemented, and $p_{\mathcal{R}}(y | a; \alpha)$ with $q(a = 1) = \alpha$ when we have a binary action set $A = \{0, 1\}$.

We could in principle assume that the agent derives beliefs from some arbitrary joint probability distribution $\hat{p}(x)$. In this case, we would have a model with exogenously fixed biased beliefs $\hat{p}(y | a)$. The personal equilibrium definition imposes restrictions on the agent’s beliefs: Through the factorization in equation (2), they must respect the agent’s strategy $q(a)$ and the extended production function. One interpretation is that the agent is experienced and thus has data on how her action impacts on the variables in her subjective model. An alternative interpretation is that there are (or have been) many other agents in the organization who exchange data with their new colleague to which she can fit her subjective model.

3 The Optimal Equilibrium Contract

In this section, we study the properties of the optimal equilibrium contract for a given extended production function $p(x_1, \dots, x_n | a)$ and subjective model \mathcal{R} . If $(w^*(y), q^*(a))$ is an optimal equilibrium contract, then $w^*(y), q^*(a)$ solve the maximization problem

$$\max_{w(y) \in W, q(a) \in \Delta(A)} \sum_{a \in A} \sum_{y \in Y} q(a) p(y | a) (y - w(y)) \quad (4)$$

subject to the constraints

$$q(a) \in \Delta(A) \text{ is a personal equilibrium at } \mathcal{R} \text{ and } w(y), \quad (IC)$$

$$\sum_{a' \in A} \sum_{y \in Y} q(a') [p_{\mathcal{R}}(y | a'; q(a)) u(w(y)) - c(a')] \geq \bar{U}. \quad (PC)$$

When the agent’s subjective model \mathcal{R} equals the objective model \mathcal{R}^* , the problem collapses to the canonical principal-agent problem, and can be solved as suggested by Grossman and Hart (1983). We first find for each pure action $a \in A$ the wage scheme $w(y)$ that implements this action at lowest possible cost. Then we choose the action-incentive scheme combination that maximizes the principal’s profit. If the agent’s subjective model \mathcal{R} differs from the objective model \mathcal{R}^* , we find the optimal equilibrium contract by applying the same procedure. However, since the agent’s beliefs $p_{\mathcal{R}}(y | a; q(a))$ may depend on the implemented strategy $q(a)$, the first step has to be done for all pure and mixed actions $q(a) \in \Delta(A)$.

Suppose the agent is risk-averse with unlimited liability, and the principal implements a

(possibly mixed) strategy $q(a)$. The Kuhn-Tucker conditions for the principal's problem are then necessary and sufficient for an optimum. Choose any action a in the support of $q(a)$. The optimal incentive scheme is then characterized by the first-order condition

$$\frac{1}{u'(w(y))} = \frac{p_{\mathcal{R}}(y; q(a))}{p(y)} \left[\mu + \sum_{a' \in A} \lambda_{a'} \frac{p_{\mathcal{R}}(y | a; q(a)) - p_{\mathcal{R}}(y | a'; q(a))}{p_{\mathcal{R}}(y; q(a))} \right] \quad (5)$$

for all $y \in Y$, where μ and $\lambda_{a'}$ are the usual Lagrange multipliers for the participation and incentive compatibility constraint, respectively. Equation (5) allows us to disentangle how a misspecification in \mathcal{R} may change the contracting problem. First, the *PC* is affected when the agent holds biased beliefs about the equilibrium distribution over output; see the first term on the right of equation (5). In Subsection 3.1, we state a sufficient condition on \mathcal{R} so that this belief is unbiased. Second, the *IC* may be affected. Suppose the principal implements a pure action a and $p_{\mathcal{R}}(y; a) = p(y)$. The ratio in the squared brackets then becomes $1 - \frac{p_{\mathcal{R}}(y|a';a)}{p_{\mathcal{R}}(y|a;a)}$, in which case the optimal incentive scheme depends on a likelihood ratio as in the canonical framework. Any difference between the contracts under the objective and subjective model is then driven by differences between the corresponding likelihood ratios. In Subsection 3.2, we examine how these differences may affect the optimal equilibrium contract.

3.1 Correct Expectations on the Equilibrium Path

We use a Bayesian network result from Spiegel (2017) that characterizes under what circumstances the agent's beliefs about the equilibrium output distribution are correct, so that $p_{\mathcal{R}}(y; q(a)) = p(y)$ for all $q(a) \in \Delta(A)$. To this end, we introduce a few definitions. A v -collider is a triple of nodes (i, j, k) such that iRj , kRj and there is no link between i and k (neither iRk nor kRi is in R). The set of v -colliders of a DAG is called its v -structure. A DAG is called perfect if it has an empty v -structure. A subset of nodes $M \subset N$ is a clique in $\mathcal{R} = (N, R)$ if iRj or jRi for any two nodes $i, j \in M$. For example, in the DAG \mathcal{R}^* from Figure 1, the set $M = \{1, 3, 4\}$ is a clique, while the set $M' = \{2, 3, 4\}$ is not. Each node is a clique in itself, so the output node n is a clique. The following result essentially restates Proposition 2 from Spiegel (2017).

Proposition 1 (Equilibrium Beliefs). *If the agent's model $\mathcal{R} = (R, N)$ is perfect, her equilibrium beliefs satisfy $p_{\mathcal{R}}(x_M; q(a)) = p(x_M)$ for all $q(a) \in \Delta(A)$ and any clique $M \subset N$.*

If the agent's subjective model \mathcal{R} is perfect, then, in a personal equilibrium, the agent correctly anticipates the marginal distribution over each variable in her model, and also the joint distribution over variables in cliques. The intuition behind this result is that perfectness ex-

cludes biased estimates due to neglect of correlation. Imagine two variables i, j that influence a third variable k . Suppose that i and j are correlated, and that the agent treats them as uncorrelated. Through the application of the factorization formula (2), the agent may then obtain a biased estimate of the marginal distribution $p(x_k)$. Perfectness implies that the agent always checks for correlations between two variables i, j when, according to her subjective model, they influence a third variable k . We obtain two useful corollaries from Proposition 1.

Corollary 1. *If the agent's model $\mathcal{R} = (R, N)$ is perfect and her equilibrium action is a pure action a^* , her equilibrium beliefs satisfy $p_{\mathcal{R}}(x_M | a^*; a^*) = p(x_M | a^*)$ for every clique $M \subset N$.*

If the equilibrium contract implements a pure strategy a^* , the agent's belief about the joint distribution of any clique M conditional on her equilibrium strategy is correct. Corollary 1 is in general not true if the equilibrium contract implements a mixed strategy $q^*(a)$. While the agent still gets the marginal equilibrium distribution over each variable right, her beliefs may also exhibit $p_{\mathcal{R}}(x_i | a'; q^*(a)) \neq p(x_i | a')$ for an action a' in the support of $q^*(a)$. Thus, the agent's expected utility conditional on a' may be biased, $\mathbb{E}_{\mathcal{R}}[u(w(y)) | a'; q^*(a)] \neq \mathbb{E}[u(w(y)) | a']$.

The second direct implication of Proposition 1 is the following result.

Corollary 2. *Suppose $(w(y), q(a))$ is an equilibrium contract. If $\mathcal{R} = (R, N)$ is perfect, the PC is satisfied at this contract if and only if this is also the case under the objective model \mathcal{R}^* .*

To see why Corollary 2 is true recall that every single node is a clique. Hence, Proposition 1 implies $p(y) = p_{\mathcal{R}}(y; q(a)) = \sum_{a \in A} q(a) p_{\mathcal{R}}(y | a; q(a))$. If \mathcal{R} is perfect, the incentive scheme therefore has to satisfy the same participation constraint as under the objective model. Thus, an agent with a misspecified – but perfect – model cannot be exploited. Throughout the paper, we will assume that \mathcal{R} is perfect. As we see next, a perfect \mathcal{R} does not imply that the principal cannot benefit from the agent's misperception.

3.2 Incentive Effects

We examine how a misspecification in the agent's subjective model \mathcal{R} can change the equilibrium contract when \mathcal{R} is perfect. By Corollary 2, only the incentive compatibility constraint can then be affected by the misspecification. We analyze a simple setting with two effort levels $a \in \{0, 1\}$, two output levels $y \in \{y_L, y_H\}$ with $y_H > y_L$, and cost $c(1) = c > c(0) = 0$. The probability of output y_H increases in the agent's effort.

Consider the marketer example from the introduction. Figure 2 shows the objective model \mathcal{R}^* and the agent's subjective model \mathcal{R} . Node 1 is the level of consumer information. It can be low ($x_1 = 0$) or high ($x_1 = 1$). Node 2 is the firm's reputation, which can be bad



Figure 2: Objective model \mathcal{R}^* (left) and subjective model \mathcal{R} (right) in the marketer example.

($x_2 = 0$) or good ($x_2 = 1$). The subjective model \mathcal{R} captures that the agent does not take reputation into account. For the objective probability distribution, we use the parametrization $p(x_i = 1 \mid x_{R(i)}) = \beta_i + \sum_{j \in R(i)} \beta_{ji} x_j$ for $i \in \{1, 2\}$ and $p(y_H \mid x_1, x_2) = \beta_3 + \beta_{13} x_1 + \beta_{23} x_2$. Making cold-calls increases consumer information, $\beta_{01} > 0$, and decreases reputation, $\beta_{02} < 0$; consumer information x_1 and reputation x_2 have a positive influence on sales, $\beta_{13} > 0$ and $\beta_{23} > 0$. We obtain the following result.

Proposition 2 (Marketer Example). *Consider the marketer example of this subsection.*

- (a) *The simplification in the agent's subjective model \mathcal{R} relaxes the IC for $\alpha = 1$.*
- (b) *The optimal equilibrium contract implements $\alpha \in \{0, 1\}$. If and only if effort costs c are small enough, the optimal equilibrium contract implements $\alpha = 1$ and the principal strictly benefits from the simplification in the agent's subjective model \mathcal{R} .*

Before we prove this result, we explain the intuition behind it and its implications. First, consider statement (a). When the principal implements $\alpha = 1$, the agent overestimates the drop in expected output when she exerts low instead of high effort. According to her subjective model \mathcal{R} , the only effect of her action on the output occurs through consumer information x_1 ; she does not take into account that a deviation to low effort would also have a positive effect on expected reputation, which translates into a positive effect on expected output. Formally, the IC under the objective model \mathcal{R}^* is

$$[\beta_{01}\beta_{13} + \beta_{02}\beta_{23}] (u(w(y_H)) - u(w(y_L))) - c \geq 0. \quad (6)$$

The term in squared brackets is the effect of effort on output and contains the consumer information channel $\beta_{01}\beta_{13}$ and the reputation channel $\beta_{02}\beta_{23}$. Under the subjective model \mathcal{R} , this second channel is missing, so that the IC becomes

$$\beta_{01}\beta_{13} (u(w(y_H)) - u(w(y_L))) \geq c. \quad (7)$$

Since the effect of effort on reputation β_{02} is negative, the simplification in \mathcal{R} relaxes the *IC*. As long as $\alpha \in (0, 1)$, the reputation effect is partly reflected in $p(y_H | x_1)$; the extent of this depends on α since α affects the correlation between consumer information and reputation. A higher correlation between consumer information and reputation would mitigate some of the effect of the agent's misperception.

Next, consider statement (b). The observation that the principal implements a pure strategy would be trivial in the canonical framework with rational expectations. This is not the case here as the agent's perceived effect of effort on output $p_{\mathcal{R}}(y_H | a = 1; \alpha) - p_{\mathcal{R}}(y_H | a = 0; \alpha)$ may vary non-monotonically in α . In the present setting, the perceived effect of effort on output is maximal at $\alpha = 1$, so that there is no reason for the principal to implement a mixed strategy. At the end of this subsection, we present an example where the unique optimal equilibrium contract indeed implements a mixed strategy $\alpha \in (0, 1)$.

Importantly, if the agent chooses a pure strategy, then, by Corollary 1 and the fact that \mathcal{R} is perfect, she correctly anticipates the joint distribution over all variables in \mathcal{R} conditional on her equilibrium action. Thus, in the data that the agent gets under the optimal equilibrium contract, there are no informational cues which could alarm her about a misspecification in her subjective model. This is a crucial difference between the present framework and models where beliefs about outcomes are biased for equilibrium actions.

Finally, the last part of statement (b) spells out that the principal strictly benefits from the simplification in \mathcal{R} when effort costs are small enough so that it is profitable to implement high effort. For a range of effort costs c , the principal implements low effort when the agent has rational expectations, but high effort if her subjective model is \mathcal{R} . This is of course not true in general. For example, if the agent's action has a positive effect on reputation, $\beta_{02} > 0$, the simplification in \mathcal{R} tightens the *IC* for $\alpha = 1$ as the agent does not take all positive effects of her action on output into account.

Proof of Proposition 2. To illustrate our approach, we present the proof of Proposition 2. We first derive $p_{\mathcal{R}}(y_H | a; \alpha)$ for a given mixed equilibrium strategy $\alpha \in (0, 1)$. The agent's equilibrium belief about the joint probability distribution of the variables in \mathcal{R} is given by $p_{\mathcal{R}}(a, x_1, y) = q(a)p(x_1 | a)p(y | x_1)$. Since node 0 and node 1 form a clique, the agent's belief about the joint probability distribution of a and x_1 is correct. Hence, $p(x_1 | a)$ is independent of α and we have $p(x_1 = 1 | a) = \beta_1 + \beta_{01}a$. However, $p(y | x_1)$ depends on α since the distribution over y also depends on x_2 . To get $p(y | x_1)$, we first derive $p(x_2 = 1 | x_1)$, i.e., the probability that $x_2 = 1$ given that value x_1 is observed at node 1 when the agent's equilibrium

action is α . We calculate

$$p(x_2 = 1 | x_1 = 1) = \frac{\alpha(\beta_1 + \beta_{01})(\beta_2 + \beta_{02}) + (1 - \alpha)\beta_1\beta_2}{\beta_1 + \alpha\beta_{01}}, \quad (8)$$

$$p(x_2 = 1 | x_1 = 0) = \frac{\alpha(1 - \beta_1 - \beta_{01})(\beta_2 + \beta_{02}) + (1 - \alpha)(1 - \beta_1)\beta_2}{1 - \beta_1 - \alpha\beta_{01}}. \quad (9)$$

With this we can calculate the equilibrium probability that output y_H realizes after observing $x_1 = 1$ and $x_1 = 0$, respectively:

$$p(y_H | x_1 = 1) = \beta_3 + \beta_{13} + \frac{\alpha(\beta_1 + \beta_{01})(\beta_2 + \beta_{02}) + (1 - \alpha)\beta_1\beta_2}{\beta_1 + \alpha\beta_{01}}\beta_{23}, \quad (10)$$

$$p(y_H | x_1 = 0) = \beta_3 + \frac{\alpha(1 - \beta_1 - \beta_{01})(\beta_2 + \beta_{02}) + (1 - \alpha)(1 - \beta_1)\beta_2}{1 - \beta_1 - \alpha\beta_{01}}\beta_{23}. \quad (11)$$

From $p_{\mathcal{R}}(a, x_1, y)$ we can now calculate the agent's subjective probability of a high output after high and low effort, respectively:

$$p_{\mathcal{R}}(y_H | a = 1; \alpha) = (\beta_1 + \beta_{01})p(y_H | x_1 = 1) + (1 - \beta_1 - \beta_{01})p(y_H | x_1 = 0), \quad (12)$$

$$p_{\mathcal{R}}(y_H | a = 0; \alpha) = \beta_1 p(y_H | x_1 = 1) + (1 - \beta_1)p(y_H | x_1 = 0). \quad (13)$$

We then use these terms to compute the *IC* for $\alpha \in (0, 1)$,

$$[p_{\mathcal{R}}(y_H | a = 1; \alpha) - p_{\mathcal{R}}(y_H | a = 0; \alpha)] (u(w(y_H)) - u(w(y_L))) = 0. \quad (14)$$

By taking the limit for $\alpha \rightarrow 1$, we obtain the *IC* for $\alpha = 1$, which is the inequality in (7). Since $\beta_{02} < 0$, this completes the proof of statement (a). To prove statement (b), note first that both *IC* and *PC* must be binding at the optimal equilibrium contract. Simple calculations show that $\beta_{01}, \beta_{13}, \beta_{23} > 0$ and $\beta_{02} < 0$ imply

$$p_{\mathcal{R}}(y_H | a = 1; \alpha) - p_{\mathcal{R}}(y_H | a = 0; \alpha) \leq \beta_{01}\beta_{13} \quad (15)$$

for all $\alpha \in (0, 1]$; that is, when the agent exerts high effort with positive probability, her perceived effect of effort on output is largest at $\alpha = 1$. The principal then cannot gain from implementing a mixed strategy. Finally, given that the optimal equilibrium contract implements either $\alpha = 0$ or $\alpha = 1$, the last part of statement (b) follows from a simple comparison of expected profits under the equilibrium contracts that implement these two actions. \square

Mixed strategy example. We show by example that it is not always optimal for the principal to implement a pure strategy. Consider again the marketer example. Assume that the agent is risk-neutral, protected by limited liability so that $w(y) \geq 0$, her outside option value is zero,

and $y_L = 0$. Suppose payoff parameters are such that the principal optimally implements some $\alpha > 0$. Standard arguments show that $w(y_L) = 0$, and that $w(y_H)$ is chosen so that the IC in (14) is satisfied. The principal's expected payoff from this contract is then

$$\mathbb{E}[V] = [\alpha p(y_H | a = 1) + (1 - \alpha)p(y_H | a = 0)] \left(y_H - \frac{c}{\Delta_{\mathcal{R}}(\alpha)} \right), \quad (16)$$

where $\Delta_{\mathcal{R}}(\alpha) = p_{\mathcal{R}}(y_H | a = 1; \alpha) - p_{\mathcal{R}}(y_H | a = 0; \alpha)$ is the agent's perceived effect of effort on output. The slope of $\Delta_{\mathcal{R}}(\alpha)$ at $\alpha = 1$ is

$$\left. \frac{d\Delta_{\mathcal{R}}(\alpha)}{d\alpha} \right|_{\alpha=1} = \beta_{01}\beta_{02}\beta_{23} \left(\frac{\beta_1}{\beta_1 + \beta_{01}} - \frac{1 - \beta_1}{1 - \beta_1 - \beta_{01}} \right). \quad (17)$$

Let the agent's action have a positive impact on both consumer information and reputation, $\beta_{01} > 0$ and $\beta_{02} > 0$. Then for $\beta_{01} \rightarrow 1 - \beta_1$ the slope in (17) converges to minus infinity. Hence, if all else equal β_{01} is sufficiently close to $1 - \beta_1$, then, starting from $\alpha = 1$, a small reduction in α reduces $w(y_H)$, and in terms of profits, this reduction overcompensates the smaller probability of high output. The optimal equilibrium contract then implements a mixed strategy. Thus, when the agent is induced to switch between periods of working hard and periods of shirking, her effort appears to her as particularly important for the final output.

Of course, when the agent chooses a mixed strategy, then the data generated in equilibrium would suffice to identify the real effect of effort on output. For this, the agent would have to analyze the data like an experimentalist and compare the average output under high and low effort, respectively. However, this "test" does not make sense according to the agent's subjective model \mathcal{R} , which does not contain a direct link between the action a and the output y . Thus, one interpretation for the mixed strategy equilibrium is that the agent does not use her data effectively to correctly derive the effect of her effort on output.

3.3 Justifiability

In our framework, the agent has a fully specified model that makes predictions about outcomes for all actions $a \in A$. A natural question is then whether the optimal equilibrium contract is also optimal for the principal when evaluated from the agent's (potentially biased) perspective. If according to her subjective beliefs the principal should have offered an alternative contract, the agent may suspect that her subjective model \mathcal{R} is not correct.⁵ We call this refinement "justifiability." It has first been defined in the unawareness literature by Filiz-Ozbay (2012). We can conveniently adapt it to our framework. In the following definition, we distinguish

⁵We do not model how in this case the agent adjusts her subjective model. One alternative is that, after becoming suspicious, she looks at the production process more closely and discovers the objective model \mathcal{R}^* .

between “justifiability” and “partial justifiability.”

Definition 3. *An equilibrium contract $(w^*(y), q^*(a))$ is justifiable at \mathcal{R} if $w^*(y), q^*(a)$ solve the maximization problem*

$$\max_{w(y) \in W, q(a) \in \Delta(A)} \sum_{a \in A} \sum_{y \in Y} q(a) p_{\mathcal{R}}(y | a; q^*(a)) (y - w(y))$$

subject to the constraints (IC) and (PC). An equilibrium contract $(w^(y), q^*(a))$ is partially justifiable at \mathcal{R} if $w^*(y)$ is a solution to this maximization problem when $q(a) = q^*(a)$ is given.*

An equilibrium contract $(w^*(y), q^*(a))$ is justifiable if the choice of incentive scheme $w^*(y)$ and implemented action $q^*(a)$ maximizes the principal’s expected payoff when evaluated according to the agent’s beliefs $p_{\mathcal{R}}(y | a; q^*(a))$. It is partially justifiable if the incentive scheme $w^*(y)$ maximizes the principal’s expected payoff, when evaluated according to the agent’s beliefs, given that the principal wants to implemented action $q^*(a)$. Partial justifiability is a weaker refinement where the agent does not start thinking about her subjective model if at least the incentive scheme appears as optimal for the principal. We examine under what circumstances an optimal equilibrium contract is (partially) justifiable, and obtain this result:

Proposition 3 (Justifiability). *Let $(w^*(y), q^*(a))$ be an optimal equilibrium contract. If we have $p_{\mathcal{R}}(y; q(a)) = p(y)$ for all $q(a) \in \Delta(A)$, the following statements hold:*

- (a) *This contract is partially justifiable at \mathcal{R} .*
- (b) *If A, Y are binary sets, $q^*(a)$ is a pure strategy, and the principal strictly prefers this contract to the optimal contract under the objective model \mathcal{R}^* , it is justifiable at \mathcal{R} .*

The proof of this result is in Appendix A.2. The first part of Proposition 3 states that an optimal equilibrium contract is partially justifiable if the agent has correct expectations on the equilibrium path. In this case, the maximization problem in (4) and that in Definition 3 are identical for any given strategy $q^*(a)$. The optimal incentive scheme that implements $w^*(y)$ then also appears to the agent as optimal for the principal. Thus, by Proposition 1, if the agent’s subjective model \mathcal{R} is perfect, the optimal equilibrium contract is partially justifiable at \mathcal{R} .

This is a significant difference to a framework where the agent’s beliefs $\hat{p}(y | a)$ are exogenously fixed. The optimal contract in such framework may not be partially justifiable since it may contain a bet that, from the agent’s perspective, is not optimal for the principal. To illustrate, consider the two-actions-two-outcomes example from the previous subsection. Suppose

that the principal implements high effort $\alpha = 1$, and that the agent's beliefs are biased so that $\hat{p}(y_H | a = 1) > p(y_H | a = 1)$ and $\hat{p}(y_H | a = 0) = p(y_H | a = 0)$. Now let effort costs c converge to zero. Under rational expectations, this would imply that the optimal contract converges to a fixed-wage contract. In contrast, under biased beliefs, the optimal contract remains bounded away from fixed wages: To exploit the agent's bias, it pays more to her after output y_H and less after output y_L . However, from the agent's perspective, an incentive scheme that is close to fixed wages would be optimal. Thus, according to her, the offered incentive scheme cannot be optimal for the principal.

To prove justifiability we additionally have to rule out that, according to the agent's beliefs, the principal can benefit by implementing a different action. Unfortunately, it is then no longer possible to derive a general statement. If an equilibrium contract is optimal for the principal, this does not imply that it is justifiable, even if the agent has correct expectations on the equilibrium path. Justifiability then has to be proven for each case individually.

The second part of Proposition 3 states sufficient (but not necessary) conditions for justifiability for a relevant special case. In a two-actions-two-outcomes setting, an optimal equilibrium contract is justifiable if it implements a pure strategy and the principal strictly benefits from the agent's misperception (which was the case in the marketer example of the previous subsection). The requirement of a pure strategy is crucial here. Consider the mixed strategy example from the previous subsection where the principal implements $\alpha \in (0, 1)$ to alter the agent's sense for the importance of her effort. From the agent's perspective, this does not make sense. According to her, it would be optimal for the principal to implement high effort with certainty. Note that she is indifferent between high and low effort, so (in her mind) the incentive scheme can remain the same. Thus, the optimal equilibrium contract in the mixed strategy example is not justifiable.

4 The Informativeness Principle

An important question in contract theory is on which information the principal should condition the agent's wage. For a setting with risk-averse agent who has unlimited liability, the informativeness principle states that the optimal contract conditions on an additional variable z if and only if it is informative about the agent's effort, i.e., if and only if the likelihood ratio $\frac{p(y,z|a')}{p(y,z|a)}$ varies in z for some y .⁶ In this section, we derive a version of the informativeness principle that allows for boundedly rational agents. To this end, we exploit the fact that an agent with biased subjective beliefs may still have correct expectations about the joint distribution of

⁶Whether this result holds or not depends on the formal details of the contracting problem; see Chaigneau et al. (2019) for a recent discussion and a further extension of the informativeness principle.

contractible variables in equilibrium. We then apply our version of the informativeness principle to provide a rationale for why in executive compensation contracts peer-performance is mostly not used so that CEOs are rewarded for windfall gains.

The original version of the informativeness principle may no longer hold when the agent's subjective model \mathcal{R} is misspecified. Consider the marketer example from Subsection 3.2 and assume that the principal can also condition the agent's wage on consumer information x_1 . If the agent had rational expectations, the optimal wage scheme would condition both on consumer information x_1 and sales x_3 since neither variable is a sufficient statistic of the other (to avoid confusion below, we here use x_3 instead of y). However, according to the agent's subjective model \mathcal{R} , sales x_3 are just a noisy signal of consumer information x_1 . Therefore, the optimal equilibrium contract only conditions on x_1 and appears as "incomplete."⁷

We can generalize this finding and obtain a version of the informativeness principle that allows for misspecified subjective models \mathcal{R} . To get this statement, we assume that the agent's subjective model is such that she correctly anticipates the joint distribution over the two contractible variables y and z . Recall from Proposition 1 that this is the case if \mathcal{R} is perfect and there is a link between y and z in \mathcal{R} (so that they form a clique).

Proposition 4 (Informativeness Principle). *Suppose the agent is risk-averse and has unlimited liability. Let y and z be two contractible variables that are both part of the agent's subjective model \mathcal{R} . If $p_{\mathcal{R}}(z, y; q(a)) = p(z, y)$ for all $q(a) \in \Delta(A)$, the following statements hold:*

- (a) *Suppose that $a \in \{0, 1\}$ and $c(1) > c(0)$. The equilibrium contract that implements $\alpha = 1$ at lowest cost to the principal does not condition on z if and only if for all triples a, y, z we have $p_{\mathcal{R}}(z | y, a; \alpha = 1) = p_{\mathcal{R}}(z | y; \alpha = 1)$.*
- (b) *If for all $q(a) \in \Delta(A)$ and all triples a, y, z we have $p_{\mathcal{R}}(z | y, a; q(a)) = p_{\mathcal{R}}(z | y; q(a))$, the optimal equilibrium contract does not condition on z .*

The proof of Proposition 4 is in Appendix A.3. We provide an interpretation of this result and explain its implications. First, the condition $p_{\mathcal{R}}(z | y, a; q(a)) = p_{\mathcal{R}}(z | y; q(a))$ for all $q(a) \in \Delta(A)$ and all triples a, y, z indicates that, in the agent's mind, variable z is independent of her action conditional on variable y (regardless of the implemented action). If this condition is satisfied, the agent believes that z does not contain any information about her action that is not already in y . However, this condition alone does not imply that the optimal equilibrium contract does not condition the agent's wage on z . In addition, the agent's subjective belief

⁷A further interesting trade-off can be observed here. Recall from the marketer example that when the contract only conditions on sales x_3 , the agent with subjective model \mathcal{R} is control optimistic, which relaxes the IC. In contrast, when the contract only conditions on consumer information x_1 , the agent has correct expectations about her expected payoff under alternative actions, so the IC is unaffected by the misspecification in \mathcal{R} .

about the joint equilibrium distribution of y and z needs to be correct. Otherwise, the principal may want to exploit the agent's biased perception of this distribution, and condition on z even if the agent thinks that z is uninformative about her action given y . This is equivalent to betting when two individuals have different prior beliefs about future events (as pointed out in the previous section, such a contract would also not be justifiable).

An interesting special case emerges when the agent believes that z is independent of all other variables.⁸ If y and z are independent in the objective model, the optimal equilibrium contract would not condition on z (even if y and z are not independent conditional on a); from the agent's perspective that would only introduce noise to the wage scheme. However, if z and y are correlated, the requirements of Proposition 4 are no longer satisfied, and the optimal equilibrium contract may imply a bet on the joint realization of y and z .

Second, Proposition 4 consists of two statements. Statement (a) is the informativeness principle for the case of binary action spaces. It is very similar to the original version: The statement implies that the optimal equilibrium contract that implements $\alpha = 1$ conditions on z if and only if the likelihood ratio $\frac{p_{\mathcal{R}}(y,z|a=0;\alpha=1)}{p_{\mathcal{R}}(y,z|a=1;\alpha=1)}$ varies in z for some y . Statement (b) for general finite action spaces is weaker since the additional information embedded in z may, according to the agent's subjective beliefs, only affect non-binding ICs.⁹

Third, observe that Proposition 4 does not impose any further assumptions on the agent's subjective model \mathcal{R} . It therefore applies to all settings in which the agent's beliefs satisfy the conditions outlined in the proposition. Importantly, we can state sufficient conditions on \mathcal{R} so that the agent's beliefs satisfy the conditional independence assumption. The Bayesian network literature establishes “ d -separation” as a convenient tool to check conditional independence of two sets of variables in a model \mathcal{R} ; we describe it in the supplementary appendix.

Fourth, our Bayesian network framework allows for a causal interpretation of the informativeness principle. The optimal equilibrium contract conditions on both y and z , if the agent's action has partially independent effects on these two variables according to \mathcal{R} ; it does not condition on z if, according to \mathcal{R} , variable z is a consequence of y . In this case, the optimal contract conditions on the variable that is “causally closer” to the agent's action.

As an application, we consider a setting in which the principal can condition the agent's wage both on her output $y \in \{y_L, y_H\}$ and on her relative performance $z \in \{-1, 0, 1\}$; the latter variable captures, for example, how the stock price of the company compares to that of the company's rivals. There is a common shock $x_1 \in \{0, 1\}$, e.g., the state of the economy, that positively affects both own output y and the rivals' output $x_3 \in \{y_L, y_H\}$. Through competition,

⁸This is the case when in her subjective model \mathcal{R} there are no links between z and any other variables in N . Note that if z were not even a variable in N , then the agent's expected utility would no longer be measurable with respect to her subjective beliefs when the contract conditions on z .

⁹This is a general issue of the informativeness principle and not specific to our framework.

output y has a negative effect on the rivals' output x_3 (e.g., if y is high, the rivals' output tends to be smaller since consumers prefer the product of the agent's firm). The objective model \mathcal{R}^* on the left in Figure 3 illustrates this setting.

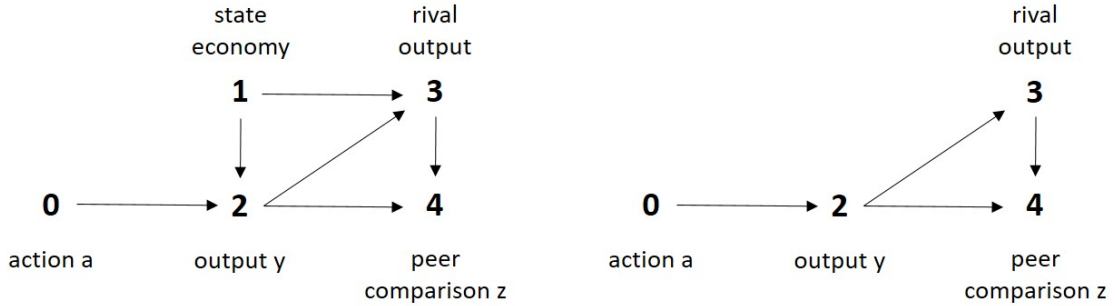


Figure 3: Objective model \mathcal{R}^* (left) and subjective model \mathcal{R} (right) in the peer-comparison example.

Under the objective model \mathcal{R}^* , the optimal equilibrium contract that implements high effort would, at any generic parametrization, condition the agent's wage both on output and relative performance. This can be established by visually inspecting \mathcal{R}^* using d -separation.¹⁰ The intuition is as follows: Suppose we know the agent's output y . Then information about the agent's action a provides additional information about the state of the economy x_1 , and hence also additional information about peer performance z . Hence, a and z are not independent conditional on y in \mathcal{R}^* .

Now suppose that the agent does not take the common shock x_1 into account so that her subjective model is given by \mathcal{R} on the right of Figure 3. Since \mathcal{R} is perfect and the variables y and z are linked in \mathcal{R} , the agent correctly anticipates the equilibrium distribution over the two variables. Moreover, if we know the output y , then, according to \mathcal{R} , the agent's action contains no further information about z (one can formally show this using d -separation). Proposition 4 then implies that the optimal equilibrium contract that implements $\alpha = 1$ only conditions on the agent's own output y . It is therefore incomplete and rewards the agent for windfall gains that come from good states of the economy. In the agent's mind, her relative performance is only a noisy signal of her own output. Hence conditioning her wage on relative performance would only increase the agent's exposure to risk and hence implementation costs.

Many actual compensation contracts indeed do not make use of peer-performance and reward executives for windfall gains. Bertrand and Mullainathan (2001) and Bebchuk and Fried (2004) discuss this phenomenon and possible explanations. A popular explanation is that executives use their influence over the board of directors to alter their compensation, which

¹⁰The "usual" way to see this is to consider a particular parametrization. Consider our linear specification with binary outcomes at all variables except z ; for z we assume that $p(z = 1 | y > x_3) \approx 1$, $p(z = 0 | y = x_3) \approx 1$, and $p(z = -1 | y < x_3) \approx 1$. If the influence of y on x_3 is small enough, the optimal contract that implements high effort conditions on both variables, and the agent's wage increases in both y and z .

then happens to increase in windfall gains. However, this theory cannot explain the inefficient risk allocation. In contrast, model misspecification can account for inefficient risk allocation. For example, the manager’s model is misspecified as in the application if she attributes the output to her action alone, or if she ignores the statistical implications of common shocks and therefore evaluates peer-performance as uninformative about her own action.

5 Behavioral Rationality

We learned in Section 3 that a simplification in the agent’s subjective model may affect the incentive compatibility constraint. However, does a simplification in \mathcal{R} automatically imply that the agent’s beliefs are biased? In this section, we show that the answer is negative. The agent may correctly anticipate the true production function even when her subjective model \mathcal{R} omits variables from \mathcal{R}^* . When this statement holds for any parametrization of the extended production function that factorizes¹¹ according to \mathcal{R}^* , we say that the agent is “behaviorally rational.” We state the formal definition.

Definition 4. *An agent with subjective model \mathcal{R} is behaviorally rational if, at any probability distribution $p(x) \in \Delta(X)$ that factorizes according to \mathcal{R}^* , we have $p_{\mathcal{R}}(y | a; q(a)) = p(y | a)$ for all $a \in A$ and $q(a) \in \Delta(A)$.*

For a given objective model \mathcal{R}^* we can characterize when the agent is behaviorally rational. We will see that two extended production functions – which involve the same set of nodes N^* and may give rise to the same production function $p(y | a)$ – can differ in the extent to which simplifications affect the agent’s beliefs about $p(y | a)$. This extent depends on the “channels” in \mathcal{R}^* through which the agent’s action affects the output. Intuitively, they describe the agent’s role in the organization, that is, which components or behaviors of others the agent affects directly or indirectly through her action. In Subsection 5.1, we motivate this interpretation in an example where the agent’s job determines the scope for biased beliefs and control optimism. In Subsection 5.2, we characterize when the agent is behaviorally rational and generalize the main findings from Subsection 5.1.

¹¹In this section, we deviate from our earlier assumption that $p(x^*)$ does not contain any additional conditional independence assumptions compared to \mathcal{R}^* . This allows us to use results and techniques from the Bayesian network literature. Importantly, if the agent is behaviorally rational in the current setting, she is also behaviorally rational under the earlier assumption.

5.1 The Agent's Job and the Scope for Control Optimism

We examine the interaction between the agent's job, model misspecification, and incentives. Let the agent first work as an ordinary marketer whose job is to increase sales. This time, making cold-calls is not part of her job. Her effort only has a (positive) effect on consumer information, for example, through informative advertising. Nevertheless, there is a group of employees engaged in telemarketing. Their effort – making cold-calls – impacts on consumer information and the firm's reputation in the usual manner. The objective model \mathcal{R}^* on the left of Figure 4a represents the causal structure of this extended production function. Throughout, we use our parametrization with binary outcomes at all variables $i \in N^*$ and $p(x_i = 1 | x_{R(i)}) = \beta_i + \sum_{j \in R(i)} \beta_{ji} x_j$. The telemarketers either conduct cold-calls or not, $\beta_1 \in \{0, 1\}$; cold-calls have a negative effect on reputation, $\beta_{13} < 0$; consumer information has a positive effect on reputation, $\beta_{23} > 0$.¹² All formal proofs of this subsection are in Appendix A.4.

Imagine that the marketer neither takes into account the telemarketers' operation nor the firm's reputation so that her subjective model is given by \mathcal{R} on the upper-left of Figure 4b. When choosing effort, she only considers the effect through consumer information. Does this misspecification change incentives? The answer is negative. We can show – using the results from the next subsection – that the agent's subjective beliefs about the production function are correct, so that $p_{\mathcal{R}}(y_H | a; \alpha) = p(y_H | a)$ for all $a \in \{0, 1\}$ and $\alpha \in [0, 1]$. Thus, given her role in the principal's project (as captured by \mathcal{R}^*), the subjective model \mathcal{R} is rich enough to produce correct predictions. The agent may ignore important parts of the project and still act as if she were fully rational. The optimal contract is then the same as in the canonical model.

Importantly, telemarketing still matters for the principal since the probability distribution over sales depends on whether cold-calls are made or not. It is just not essential for the agent to know whether cold-calls take place. Her estimate of the production function implicitly takes into account the deterministic activity of the telemarketers.

Is there any simplification that would make the agent overestimate the effectiveness of her effort? Again, the answer is negative. If the agent does not take node 2 into account, she believes that her action has no consequences for the output. It would then be impossible to implement high effort. If only node 1 or only node 3 were omitted from her subjective model, the agent would again have correct beliefs about the production function. Thus, there is no scope for control optimism when the agent works as ordinary marketer.

¹²Here we introduce a link between consumer information and reputation, and violate our full support assumption by assuming $p(x_1 = 1) \in \{0, 1\}$. The latter implies that in objective model \mathcal{R}^* we could drop node 1 and factor the value $p(x_1 = 1)$ into the other conditional probabilities. If $p(x_1 = 1) \in (0, 1)$, node 1 would be a confounding factor and the behavioral rationality result in the example would no longer hold (we consider such a case in the second application in Section 6). In terms of interpretation, this assumption just means telemarketing either takes place or not.

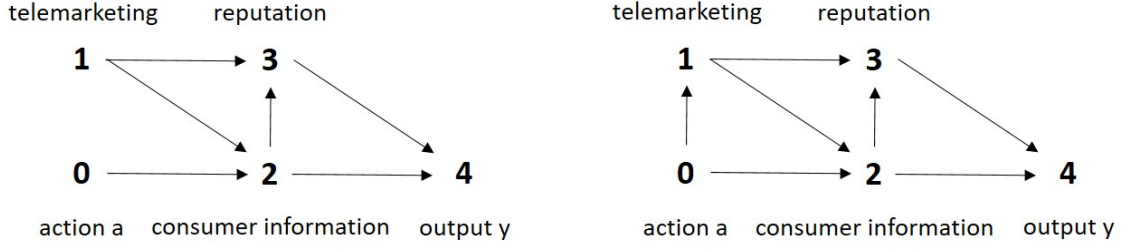


Figure 4a: Objective model \mathcal{R}^* (left) when the agent works as ordinary marketer, and objective model \mathcal{R}^{**} (right) when the agent works as “head of marketing.”

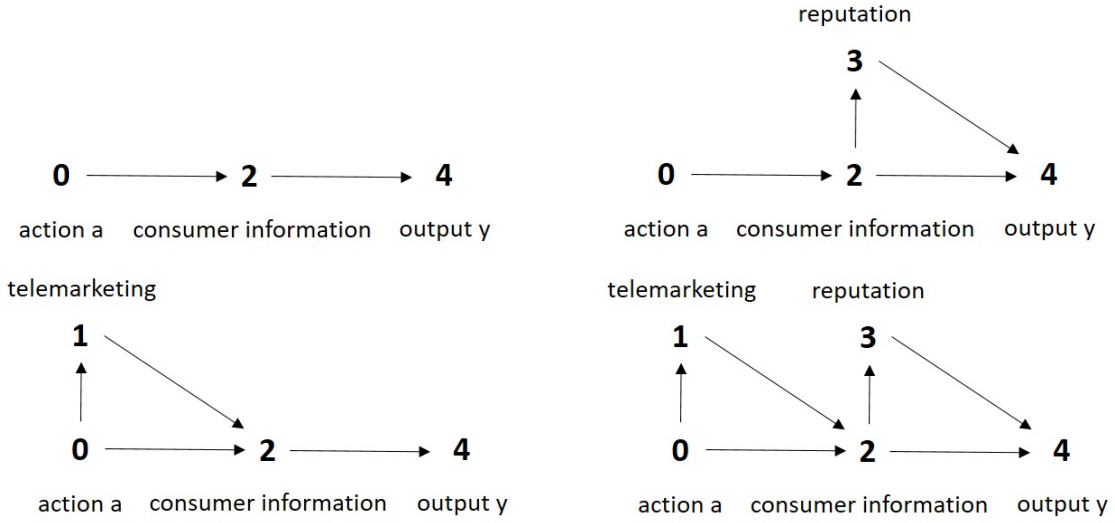


Figure 4b: Subjective models \mathcal{R} (upper-left), \mathcal{R}_1 (upper-right), \mathcal{R}_2 (lower-left), and \mathcal{R}_3 (lower-right).

Next, we alter the agent’s job by promoting her to “head of marketing.” Her action now influences the telemarketers’ effort, for example, by motivating or inspiring the telemarketers. Instead of $p(x_1 = 1) = \beta_1$, we now have $p(x_1 = 1 | a) = \beta_1 + \beta_{01}a$. To keep things as close as possible to the previous case, we assume $\beta_1 = 0$ and $\beta_{01} = 1$.¹³ Hence, the agent needs to act in order to get the telemarketers going. The objective model of the extended production function is given by \mathcal{R}^{**} on the right of Figure 4a. How does a misspecification in the agent’s subjective model now affect equilibrium beliefs and incentives in this environment?

Let us first assume that the agent has the same subjective model \mathcal{R} as before (on the upper-left of Figure 4b). She neglects both the telemarketers’ activity and the firm’s reputation. This is not realistic since as “head of marketing” the agent should be aware of her subordinates’ basic activities; so we will relax this assumption below. The misspecification now affects incentives. Under the objective model \mathcal{R}^{**} the IC that implements $\alpha = 1$ would be

$$[(\beta_{02} + \beta_{01}\beta_{12})(\beta_{24} + \beta_{23}\beta_{34}) + \beta_{01}\beta_{13}\beta_{34}][u(w(y_H)) - u(w(y_L))] \geq c. \quad (18)$$

¹³Formally, we assume $\beta_1 = \varepsilon_1$ and $\beta_{01} = 1 - \varepsilon_2$ where $\varepsilon_1 < \varepsilon_2$, and consider the limit beliefs as $\varepsilon_1 \rightarrow 0$ and $\varepsilon_2 \rightarrow 0$. We show in the proofs for this subsection that our results do not depend on this assumption.

The squared brackets contain the different channels through which effort affects output. The partial negative effect of effort on output through cold-calls and reputation is captured in the term $\beta_{01}\beta_{13}\beta_{34}$; it is negative since $\beta_{13} < 0$. Under the subjective model \mathcal{R} the *IC* becomes

$$(\beta_{02} + \beta_{01}\beta_{12})(\beta_{24} + \beta_{23}\beta_{34})(u(w(y_H)) - u(w(y_L))) \geq c. \quad (19)$$

Here the partial negative effect is missing so that the *IC* is relaxed. Note that through the estimate of the link between the agent's action and consumer information, the agent implicitly takes into account her positive influence on the telemarketers' effort, which in turn positively affects consumer information (see the term $\beta_{01}\beta_{12}$). Therefore, by being promoted to a job where the agent also influences telemarketing, she overestimates her productivity. The principal benefits from this since the misspecification reduces the need to provide effort incentives.

Assume now that the agent takes the telemarketers' action into account, but still omits reputation in her model. Therefore, her subjective model is given by \mathcal{R}_2 on the lower-left of Figure 4b. Does this inclusion correct, at least partly, the agent's beliefs? It turns out that this is not the case. The models \mathcal{R} and \mathcal{R}_2 produce the same beliefs about the effectiveness of effort, i.e., $p_{\mathcal{R}}(y_H | a; \alpha) = p_{\mathcal{R}_2}(y_H | a; \alpha)$ for all $a \in \{0, 1\}$ and $\alpha \in [0, 1]$. Including more variables does not necessarily make the agent more rational. This also holds for the models \mathcal{R}_1 and \mathcal{R}_3 in Figure 4b. Note that \mathcal{R}_3 is almost equal to the objective model \mathcal{R}^{**} , only the link between telemarketing and reputation is missing. Yet, all subjective models in this figure produce the same beliefs. Thus, a small misspecification in the agent's subjective model can render several important variables as inessential for estimating the production function.

Proposition 5 (Scope for Control Optimism). *Consider the job examples of this subsection.*

- (a) *If the agent works as ordinary marketer (objective model \mathcal{R}^*), the misspecification in \mathcal{R} has no effect on the *IC* and the optimal equilibrium contract is the same as in the canonical model. There is no simplification that generates control optimism.*
- (b) *If the agent works as "head of marketing" (objective model \mathcal{R}^{**}), the misspecification in \mathcal{R} generates control optimism and relaxes the *IC*; the subjective models \mathcal{R} , \mathcal{R}_1 , \mathcal{R}_2 , and \mathcal{R}_3 generate the same beliefs about the production function.*

Proposition 5 illustrates how the agent's job may matter for optimal incentives. The two jobs with objective models \mathcal{R}^* and \mathcal{R}^{**} may give rise to the same production function $p(y | a)$,¹⁴ so that incentives would be identical under rational expectations. However, effort motivation is

¹⁴Specifically, when we denote parameters for the job with objective model \mathcal{R}^* (\mathcal{R}^{**}) with "****" ("****") we only have to select parameters so that $\beta_{02}^*(\beta_{24}^* + \beta_{23}^*\beta_{34}^*) = (\beta_{02}^{**} + \beta_{01}^{**}\beta_{12}^{**})(\beta_{24}^{**} + \beta_{23}^{**}\beta_{34}^{**}) + \beta_{01}^{**}\beta_{13}^{**}\beta_{34}^{**}$.

larger under a job with the objective model \mathcal{R}^{**} when the agent's subjective model is simplified in a way that benefits the principal. The crucial difference between the jobs are the sets of channels through which the action affects the output. In the next subsection, we will formally define these channels.

Part (a) and (b) of Proposition 5 combined demonstrate that an agent's degree of control optimism may be determined by the nature of her job. In the example, the agent with misspecified model \mathcal{R} was behaviorally rational in her job as ordinary marketer, but overestimated the importance of her effort after being promoted to "head of marketing" where she influences the actions of others. Thus, in our framework, the agent's control optimism is not caused by certain features of her personality, but it is a consequence of her environment when her subjective model does not capture all empirical regularities of this environment.

5.2 A General Result on Behavioral Rationality

To obtain a general result on behavioral rationality, we assume that the objective model \mathcal{R}^* is perfect, and that the agent's subjective model \mathcal{R} is a simplification. \mathcal{R} will then be perfect. No v -structure emerges if we take out nodes from a perfect \mathcal{R}^* and all links attached to them. The assumptions on \mathcal{R}^* and \mathcal{R} are not overly restrictive: Any probability distribution $p(x^*)$ factorizes according to some perfect DAG \mathcal{R}^* . The assumption on \mathcal{R} is satisfied by almost all subjective models in this paper. We can also (partially) extend our behavioral rationality result to imperfect objective models. All formal proofs for this subsection are in Appendix A.5.

In the following, we characterize for any perfect \mathcal{R}^* the subset of nodes the agent needs to have in her subjective model \mathcal{R} so that she acts as if she had fully rational beliefs about the production function. We use the following definitions and results from the Bayesian network literature. Consider any DAG $\mathcal{R} = (N, R)$. Its skeleton (N, \tilde{R}) is obtained by making the DAG undirected. We have $i\tilde{R}j$ if and only if iRj or jRi .

Definition 5. *Two DAGs \mathcal{R} and \mathcal{G} are equivalent if $p_{\mathcal{R}}(x) \equiv p_{\mathcal{G}}(x)$ for every $p(x) \in \Delta(X)$.*

Proposition 6 (Verma and Pearl 1991). *Two DAGs \mathcal{R} and \mathcal{G} are equivalent if and only if they have the same skeleton and v -structure.*

Two different models produce the same beliefs if they share the same skeleton and the same set of v -colliders. A subset of nodes $M \subset N$ is called ancestral in \mathcal{R} if for all nodes $i \in M$ we have $R(i) \subset M$. A path τ of length d from node i to node j is a sequence of nodes $\tau_0, \tau_1, \dots, \tau_d$ so that $\tau_0 = i$, $\tau_d = j$, and $\tau_{h-1}\tilde{R}\tau_h$ for all $h \in \{1, \dots, d\}$. The length of the shortest path between i and j is called the distance between these nodes and denoted by $d(i, j)$. A path of length d is active if there is no $h \in \{1, \dots, d-1\}$ so that $\tau_{h-1}R\tau_h$ and $\tau_{h+1}R\tau_h$.

Define by \mathcal{E} the set of DAGs in the equivalence class of \mathcal{R}^* in which the action node 0 is ancestral (nothing influences the agent's action). In each of these DAGs, all active paths between the action node 0 and any node i point towards i . Thus, the assumption that node 0 is ancestral pins down the direction of many links in a perfect DAG. We call such links “fundamental links.” There is a close connection between fundamental links and the set of nodes that can be removed while maintaining behavioral rationality.

Definition 6. *Consider two nodes $i, j \in N^*$. If iGj for all $\mathcal{G} = (G, N^*) \in \mathcal{E}$, then the link iGj is called fundamental link and denoted by iEj .*

An intuition for fundamental links is that they capture empirically relevant directions of causality (given agreement on the ancestral node). Specifically, they describe how the agent's action impacts on other variables. Consider \mathcal{R}^* from Figure 1. Since the action node is ancestral, the links pointing from node 0 to other nodes are fundamental ($0R^*1$, $0R^*2$, and $0R^*3$). Thus, the two links pointing into the output node ($1R^*4$ and $3R^*4$) also must be fundamental. If we would turn around one of them, we would create a v -collider since there is no link between node 0 and node 4. The remaining links $1R^*2$, $1R^*3$, and $2R^*3$ are not fundamental. We can state a result that characterizes all fundamental links in any perfect DAG; see Appendix A.5. For now, we go a step further and consider sequences of fundamental links.

Definition 7. *Let τ be an active path in \mathcal{R}^* . Then τ is a fundamental active path if all the links between neighboring nodes in τ are fundamental.*

Fundamental active paths are what we so far called “channels.” Consider again \mathcal{R}^* from Figure 1. The path $\tau = \{0, 1, 4\}$ is a fundamental active path since both links $0R^*1$ and $1R^*4$ are fundamental. In contrast, the active path $\tau' = \{0, 2, 3, 4\}$ is not fundamental since the link $2R^*3$ is not fundamental. We define the set of nodes that are part of at least one fundamental active path between the action and the output by

$$H^*(\mathcal{R}^*) := \{i \in N^* \mid i \text{ is part of a fundamental active path between } 0 \text{ and } n \text{ in } \mathcal{R}^*\}.$$

It turns out that the nodes in $H^*(\mathcal{R}^*)$ are exactly those nodes the agent needs to have in her subjective model in order to be behaviorally rational, provided that her subjective model is a simplification. We can prove this by finding a DAG \mathcal{G} that is equivalent to \mathcal{R}^* and in which there are no links pointing from nodes in $N^* \setminus H^*(\mathcal{R}^*)$ to nodes in $H^*(\mathcal{R}^*)$. In this DAG, the nodes that are not in $H^*(\mathcal{R}^*)$ have no influence on the output, so the agent can safely ignore them. By Proposition 6, the agent correctly anticipates the production function if $H^*(\mathcal{R}^*) \subseteq N$.

Proposition 7 (Behavioral Rationality). *Let \mathcal{R}^* be a perfect DAG and let the agent's subjective DAG \mathcal{R} be a simplification. The agent is behaviorally rational if and only if \mathcal{R} contains all nodes from $H^*(\mathcal{R}^*)$.*

Proposition 7 implies that the agent does not necessarily have to take into account all variables of her (potentially) complex environment in order to be behaviorally rational. In particular, this holds independent of the parametrization of the extended production function. For example, when $p(x_1, \dots, x_4 | a)$ factorizes according to \mathcal{R}^* in Figure 1, the agent can ignore node 2 and still would behave as in the contracting model with common priors. The intuition is that when $H^*(\mathcal{R}^*) \subseteq N$, then the information captured through the variables in $H^*(\mathcal{R}^*)$ already includes the probabilistic information from variables outside $H^*(\mathcal{R}^*)$. Conversely, if the agent's subjective model does not include all variables from $H^*(\mathcal{R}^*)$, she is not behaviorally rational. In this case, we can find a parametrization of $p(x_1, \dots, x_n | a)$ such that the incentive compatibility constraint is affected by the simplification in the agent's subjective model \mathcal{R} .

Next, Proposition 7 also shows that different misspecifications can have the same effect on incentives. Consider the two models \mathcal{R}_1 and \mathcal{R}_2 from the job example in Figure 4b. The set of nodes on fundamental active paths is the same for these two models, $H^*(\mathcal{R}_1) = H^*(\mathcal{R}_2) = \{0, 2, 4\}$. This implies that the agent's beliefs under these models are identical. Thus, it does not matter for the equilibrium contract whether the agent ignores node 1, node 3, or both nodes. Therefore, the ignorance about one channel of causality may render another variable unimportant. A further interpretation is that two agents with different subjective models may have the same beliefs about the production function. We capture this result in a general statement. Consider a DAG $\mathcal{R} = (N, R)$ and a subset $\tilde{N} \subset N$. Denote by $\mathcal{R}^{[\tilde{N}]} = (\tilde{N}, \tilde{R})$ with $\tilde{R} = (\tilde{N} \times \tilde{N}) \cap R$ the DAG \mathcal{R} restricted on \tilde{N} .

Corollary 3. *Let $\mathcal{R}_1 = (N_1, R_1)$ and $\mathcal{R}_2 = (N_2, R_2)$ be two perfect DAGs. Suppose there exists a DAG \mathcal{R}_3 so that $\mathcal{R}_3^{[N_1]} = \mathcal{R}_1$ and $\mathcal{R}_3^{[N_2]} = \mathcal{R}_2$. If $H^*(\mathcal{R}_1) = H^*(\mathcal{R}_2)$, then we have that $p_{\mathcal{R}_1}(y | a; q(a)) = p_{\mathcal{R}_2}(y | a; q(a))$ for all $a \in A$ and $q(a) \in \Delta(A)$.*

Finally, note that one can make any imperfect DAG perfect by adding links between nodes that create v -colliders. We can exploit this to partially extend Proposition 7 to imperfect objective models; see the supplementary appendix for an elaborate discussion.

6 Comparative Statics

One advantage of our approach to contracting with boundedly rational agents is that beliefs are derived endogenously from the true production process. This allows us to analyze how the

optimal equilibrium contract varies in the parameters of the environment. In this section, we revisit two comparative statics that have received considerable attention in the literature: the trade-off between risk and incentives, and the relationship between team size and incentives. In both cases, the empirical evidence on these comparative statics conflicts with the predictions of the canonical model. We briefly discuss how we can explain these findings within our framework. All formal details of this section are relegated to the supplementary appendix.

Risk and Incentives. A risk-averse agent demands a risk premium for accepting a wage schedule with uncertain wage payments. Thus, an increase in risk drives up the costs of providing incentives. Consequently, the provision of effort incentives should decrease in the riskiness of the environment. However, empirically this relationship does not hold in general (e.g., Prendergast 2002). Field evidence on the relationship between risk and incentives for CEO compensation is mixed, and for other domains, such as franchising, a positive relationship can be observed. In contrast, a negative relationship is obtained in lab experiments where subjects know the true production function (Corgnet and Hernán-González 2019).

We can use our marketer example to show how the relationship between risk and incentives may become positive when the agent has a simplified model of the project. We consider a mean-preserving spread in $p(y | a)$, so that under the objective model \mathcal{R}^* the provision of incentives becomes more costly when there is more risk. However, if the agent's subjective model is misspecified, there can be an additional effect of risk on incentives: The agent may perceive the riskier environment as one in which her action is more important for the output, which relaxes the incentive compatibility constraint. If this effect is sufficiently strong relative to the risk premium effect, there can be a positive relationship between risk and incentives.

Team Size and Incentives. In a team incentive problem, effort incentives are provided by tying each team member's payoff to the joint output y . The effectiveness of team incentives is constrained by the size of the team. When an agent's relative contribution to the output becomes small, it is typically no longer profitable for the principal to condition her pay on y , as the incentive effect would be outweighed by the costs of incentive provision (e.g., Kandel and Lazear 1992). An important implication of this result is that stock-options should be granted only to those employees whose actions significantly move the stock price. However, many firms grant stock options also to non-executive employees, and there is evidence that these have positive incentive effects (e.g., Hochberg and Lindsey 2010).

We can provide a belief-based explanation for this phenomenon in a setting with many agents. Each agent produces an intermediate output which positively affects the final output. A common shock affects all intermediate outputs in the same direction. If an agent ignores the intermediate outputs by other agents, she perceives a strong relationship between her interme-

diate output and the final output. She then overestimates the importance of her effort, which relaxes the incentive compatibility constraint. We demonstrate that output-based incentives then can remain effective even when the team becomes arbitrary large.

7 Conclusion

In this paper, we applied Spiegler's (2016) Bayesian network framework to analyze optimal contracting in a principal-agent setting where the agent forms beliefs about the production function based on a misspecified model of the principal's project. The objective causal model may be very complex, and may contain empirical regularities that the agent does not consider due to cognitive limitations or because they are never brought to her attention.

The optimal contract exhibits the following features. First, it does not exploit the agent if her subjective model takes into account the correlation between variables in her model that have a joint influence on a third variable (in which case it is "perfect"). Second, the principal may nevertheless benefit from a misspecification in the agent's perfect subjective model if it makes the agent control optimistic so that the incentive compatibility constraint is relaxed. Third, if the agent's subjective model is perfect, the agent cannot infer from the shape of incentives that her beliefs are biased. Fourth, when the agent correctly anticipates the joint distribution of contractible variables, the optimal contract conditions on an additional variable only if it is informative about the action according to the agent's model. Fifth, the optimal contract is identical to the rational benchmark if the agent is behaviorally rational. We characterize when this is the case, and apply this finding to show how the scope for control optimism may depend on the agent's job. For example, a front-line worker may not fully understand the workings of the organization around her, but still act as if she were fully rational. In contrast, a high-ranking manager, who affects the output by influencing the behavior of many subordinates, overestimates her own productivity if she does not take into account the challenges that her subordinates face in their routines.

We focused on a simple contracting framework so that we can identify precisely how misspecifications in the agent's model affect incentive contracts. Future research can extend the framework by considering team incentives, relational contracts, and delegation. The Bayesian network approach offers a very disciplined tool to study the effects of bounded rationality on organizations, and we think that our results are useful in this respect.

References

- AUSTER, SARAH (2013), “Asymmetric awareness and moral hazard.” *Games and Economic Behavior*, 82, 503–521.
- BEBCHUK, LUCIAN ARYE, AND JESSE FRIED (2004), *Pay without performance: The unfulfilled promise of executive compensation*. Harvard University Press, Cambridge.
- BÉNABOU, ROLAND (2013), “Groupthink: Collective delusions in organizations and markets.” *Review of Economic Studies*, 80, 429–462.
- BÉNABOU, ROLAND, AND JEAN TIROLE (2002), “Self-confidence and personal motivation.” *Quarterly Journal of Economics*, 117, 871–915.
- BERTRAND, MARIANNE, AND SENDHIL MULLAINATHAN (2001), “Are CEOs rewarded for luck? The ones without principals are.” *Quarterly Journal of Economics*, 116, 901–932.
- BLAKE, THOMAS, CHRIS NOSKO, AND STEVEN TADELIS (2015), “Consumer heterogeneity and paid search effectiveness: A large-scale field experiment.” *Econometrica*, 83, 155–174.
- BLOOM, NICHOLAS, BENN EIFERT, APRAJIT MAHAJAN, DAVID MCKENZIE, AND JOHN ROBERTS (2013), “Does management matter? Evidence from India.” *Quarterly Journal of Economics*, 128, 1–51.
- BRUNNERMEIER, MARKUS, AND JONATHAN PARKER (2005), “Optimal expectations.” *American Economic Review*, 95, 1092–1118.
- CHAIGNEAU, PIERRE, ALEX EDMANS, AND DANIEL GOTTLIEB (2019), “The informativeness principle without the first-order approach.” *Games and Economic Behavior*, 113, 743–755.
- CORNET, BRICE, AND ROBERTO HERNÁN-GONZÁLEZ (2019), “Revisiting the trade-off between risk and incentives: The shocking effect of random shocks?” *Management Science*, 65, 1096–1114.
- DE LA ROSA, ENRIQUE (2011), “Overconfidence and moral hazard.” *Games and Economic Behavior*, 73, 429–451.
- DEKEL, EDDIE, BARTON LIPMAN, AND ALDO RUSTICHINI (1998), “Standard state-space models preclude unawareness.” *Econometrica*, 66, 159–173.
- ELIAZ, KFIR, AND RANI SPIEGLER (2020), “A model of competing narratives.” *American Economic Review*, forthcoming.

- ELIAZ, KFIR, RANI SPIEGLER, AND HEIDI THYSEN (2019), “Strategic interpretations.” CEPR Discussion Paper No. 13441.
- FANG, HANMING, AND GIUSEPPE MOSCARINI (2005), “Morale hazard.” *Journal of Monetary Economics*, 52, 749–777.
- FILIZ-OZBAY, EMEL (2012), “Incorporating unawareness into contract theory.” *Games and Economic Behavior*, 76, 181–194.
- GERVAIS, SIMON, AND ITAY GOLDSTEIN (2007), “The positive effects of biased self-perceptions in firms.” *Review of Finance*, 11, 453–496.
- GROSSMAN, SANFORD, AND OLIVER HART (1983), “An analysis of the principal-agent problem.” *Econometrica*, 51, 7–45.
- HANNA, REMA, SENDHIL MULLAINATHAN, JOSHUA SCHWARTZSTEIN (2014), “Learning through noticing: Theory and evidence from a field experiment.” *Quarterly Journal of Economics*, 129, 1311–1353.
- HEIFETZ, AVIAD, MARTIN MEIER, AND BURKHARD SCHIPPER (2006), “Interactive unawareness.” *Journal of Economic Theory*, 130, 78–94.
- HEIFETZ, AVIAD, MARTIN MEIER, AND BURKHARD SCHIPPER (2013), “Unawareness, beliefs, and speculative trade.” *Games and Economic Behavior*, 77, 100–121.
- HOCHBERG, YAEL, AND LAURA LINDSEY (2010), “Incentives, targeting, and firm performance: An analysis of non-executive stock options.” *Review of Financial Studies*, 23, 4148–4186.
- HOLMSTRÖM, BENGT (1979), “Moral hazard and observability.” *Bell Journal of Economics*, 10, 74–91.
- IMMORDINO, GIOVANNI, ANNA MARIA MENICHINI, MARIA GRAZIA ROMANO (2015), “Contracts with wishful thinkers.” *Journal of Economics and Management Strategy*, 24, 863–886.
- KANDEL, EUGENE, AND EDWARD LAZEAR (1992), “Peer pressure and partnerships.” *Journal of Political Economy*, 100, 801–817.
- KOMINERS, SCOTT (2017), “Uber really wants you to use its credit card.” Published on bloomberg.com, downloaded on October 7 from www.bloomberg.com/opinion/articles/2017-12-06/uber-really-wants-you-to-use-its-credit-card.

- KŐSZEGI, BOTOND (2006), “Ego utility, overconfidence, and task choice.” *Journal of the European Economic Association*, 4, 673–707.
- KŐSZEGI, BOTOND (2014), “Behavioral contract theory.” *Journal of Economic Literature*, 52, 1075–1118.
- MIKLÓS-THAL, JEANINE, AND JUANJUAN ZHANG (2013), “(De)marketing to manage consumer quality inferences.” *Journal of Marketing Research*, 50, 55–69.
- NULAND, SHERWIN (2004), *The doctors’ plague: Germs, childbed Fever, and the strange story of Ignac Semmelweis*. W. W. Norton Company.
- PEARL, JUDEA (2009), *Causality: Models, Reasoning, and Inference*. Cambridge University Press.
- PRENDERGAST, CANICE (2002), “The tenuous trade-off between risk and incentives.” *Journal of Political Economy*, 110, 1071–1102.
- SANTOS-PINTO, LUÍS (2008), “Positive self-image and incentives in organisations.” *Economic Journal*, 118, 1315–1332.
- SAUTMANN, ANJA (2007), “Self-confidence in a principal-agent relationship.” Unpublished manuscript, Brown University.
- SAUTMANN, ANJA (2013), “Contracts for agents with biased beliefs: Some theory and an experiment.” *American Economic Journal: Microeconomics*, 5, 124–156.
- SCHENONE, PABLO (2020), “Causality: A decision theoretic framework.” Working Paper, California Institute of Technology.
- SIMON, HERBERT (1947), *Administrative Behavior*. Macmillan, London.
- SIMON, HERBERT (1955), “A behavioral model of rational choice.” *Quarterly Journal of Economics*, 69, 99–118.
- SPIEGLER, RAN (2016), “Bayesian networks and boundedly rational expectations.” *Quarterly Journal of Economics*, 131, 1243–1290.
- SPIEGLER, RAN (2017), “Data monkeys: A procedural model of extrapolation from partial statistics.” *Review of Economic Studies*, 84, 1818–1841.
- SPIEGLER, RAN (2019), “Can agents with causal misperceptions be systematically fooled?” *Journal of the European Economic Association*, 18, 583–617.

SPINNEWIJN, JOHANNES (2013), “Insurance and perceptions: How to screen optimists and pessimists.” *Economic Journal*, 123, 606–633.

SPINNEWIJN, JOHANNES (2015), “Unemployed but optimistic: Optimal insurance design with biased beliefs.” *Journal of the European Economic Association*, 13, 130–167.

VON THADDEN, ERNST-LUDWIG, AND XIAOJIAN ZHAO (2012), “Incentives for unaware agents.” *Review of Economic Studies*, 79, 1151–1174.

VON THADDEN, ERNST-LUDWIG, AND XIAOJIAN ZHAO (2014), “Multitask agency with unawareness.” *Theory and Decision*, 77, 197–222.

VAN DEN STEEN, ERIC (2005), “Organizational beliefs and managerial vision.” *Journal of Law, Economics, and Organization*, 21, 256–283.

VERMA, THOMAS, AND JUDEA PEARL (1991), “Equivalence and synthesis of causal models.” *Uncertainty in Artificial Intelligence*, 6, 255–268.

A Appendix

A.1 Existence of a Personal Equilibrium

We show that a personal equilibrium exists at any admissible \mathcal{R} and $w(y) \in W$. Note that $\Delta(A)$ is non-empty, compact, and convex. Define the best-response correspondence $BR : \Delta(A) \rightarrow \Delta(A)$ by

$$BR(q(a)) = \arg \max_{\tilde{q}(a') \in \Delta(A)} \sum_{a' \in A} \sum_{y \in Y} \tilde{q}(a') [p_{\mathcal{R}}(y | a'; q(a)) u(w(y)) - c(a')]. \quad (\text{A.1})$$

For every $q(a) \in \Delta(A)$ we have that $BR(q(a))$ is non-empty and convex. The latter statement follows since any convex combination of pure actions that are optimal for the agent is an element of $BR(q(a))$. Definition 1 and the factorization formula in (2) imply that the agent's beliefs $p_{\mathcal{R}}(y | a'; q(a))$ are continuous in $q(a)$. Therefore, we also must have that $\sum_{a' \in A} \sum_{y \in Y} \tilde{q}(a') [p_{\mathcal{R}}(y | a'; q(a)) u(w(y)) - c(a')]$ is continuous in $q(a)$. Hence, $BR(q(a))$ is upper hemi-continuous. The existence of a personal equilibrium then follows from Kakutani's theorem.

A.2 Omitted Proofs from Section 3

Proof of Proposition 3. Statement (a) is proven in the main text. We prove statement (b). We denote $A = \{0, 1\}$ and $Y = \{0, 1\}$ with the usual interpretation. Since the principal strictly prefers $(w^*(y), p^*(a))$ to the optimal contract under the objective model \mathcal{R}^* , and the agent correctly anticipates the equilibrium distribution over output, the equilibrium action must be $a^* = 1$ and $w^*(1) > w^*(0)$. We show that from the agent's perspective the principal cannot gain by implementing $a = 0$. Denote by \bar{w} the fixed wage that implements $a = 0$ at lowest costs to the principal under the objective model. The agent anticipates that a fixed wage of \bar{w} would optimally implement $a = 0$. Since Y is binary we must have $p(y = 1 | a = 0) > p_{\mathcal{R}}(y = 1 | a = 0; a^*)$. Thus, we get

$$\begin{aligned} \sum_{y \in Y} p_{\mathcal{R}}(y | a = 1; a^*)(y - w^*(y)) &= \sum_{y \in Y} p(y | a = 1)(y - w^*(y)) \\ &> \sum_{y \in Y} p(y | a = 0)(y - \bar{w}) \\ &> \sum_{y \in Y} p_{\mathcal{R}}(y | a = 0; a^*)(y - \bar{w}), \end{aligned}$$

where the first inequality follows from the fact that the principal strictly prefers $(w^*(y), p^*(a))$ to the optimal contract under model \mathcal{R}^* . This completes the proof of statement (b). \square

A.3 Omitted Proofs from Section 4

Proof of Proposition 4. We first prove statement (b). Suppose the principal wishes to implement $q(a)$. Since the agent is risk-averse with unlimited liability and her action set A is finite, we can use the arguments in Grossman and Hart (1983) to show that the Kuhn-Tucker theorem yields necessary and sufficient conditions for an optimum. The optimal incentive scheme is therefore characterized by the first-order condition

$$\frac{1}{u'(w(y, z))} = \frac{p_{\mathcal{R}}(y, z; q(a))}{p(y, z)} \left[\mu + \sum_{a' \in A} \lambda_{a'} \frac{p_{\mathcal{R}}(y, z | a; q(a)) - p_{\mathcal{R}}(y, z | a'; q(a))}{p_{\mathcal{R}}(y, z; q(a))} \right]. \quad (\text{A.2})$$

By assumption, we have $p_{\mathcal{R}}(y, z; q(a)) = p(y, z)$. We can rewrite $p_{\mathcal{R}}(y, z | a; q(a))$ as

$$p_{\mathcal{R}}(y, z | a; q(a)) = p_{\mathcal{R}}(y | a; q(a))p_{\mathcal{R}}(z | y, a; q(a)) = p_{\mathcal{R}}(y | a; q(a))p_{\mathcal{R}}(z | y; q(a)), \quad (\text{A.3})$$

where the last equality follows from the assumption $p_{\mathcal{R}}(z | y, a; q(a)) = p_{\mathcal{R}}(z | y; q(a))$ for all triples a, y, z . Similarly, we can write $p_{\mathcal{R}}(y, z; q(a)) = p_{\mathcal{R}}(y; q(a))p_{\mathcal{R}}(z | y; q(a))$. Hence, we get

$$p_{\mathcal{R}}(y, z | a; q(a)) - p_{\mathcal{R}}(y, z | a'; q(a)) = \frac{p_{\mathcal{R}}(y, z; q(a))}{p_{\mathcal{R}}(y; q(a))} [p_{\mathcal{R}}(y | a; q(a)) - p_{\mathcal{R}}(y | a'; q(a))]. \quad (\text{A.4})$$

The first-order condition in (A.2) therefore simplifies to

$$\frac{1}{u'(w(y, z))} = \mu + \sum_{a' \in A} \lambda_{a'} \frac{p_{\mathcal{R}}(y | a; q(a)) - p_{\mathcal{R}}(y | a'; q(a))}{p_{\mathcal{R}}(y; q(a))}. \quad (\text{A.5})$$

Since the right-hand side of this first-order equation is independent of z , the optimal incentive scheme does not condition on z , which completes the proof. Next, we prove statement (a). Risk-aversion and unlimited liability imply that the optimal incentive scheme that implements $a = 1$ is characterized by the first-order condition

$$\frac{1}{u'(w(y, z))} = \frac{p_{\mathcal{R}}(y, z | a = 1; \alpha = 1)}{p(y, z | a = 1)} \left[\mu + \lambda \left(1 - \frac{p_{\mathcal{R}}(y, z | a = 0; \alpha = 1)}{p_{\mathcal{R}}(y, z | a = 1; \alpha = 1)} \right) \right], \quad (\text{A.6})$$

where μ, λ are strictly positive constants. As above, we can write $p_{\mathcal{R}}(y, z | a = 1; \alpha = 1) = p(y, z | a = 1)$, so that this first-order condition simplifies to

$$\frac{1}{u'(w(y, z))} = \mu + \lambda \left(1 - \frac{p_{\mathcal{R}}(y, z | a = 0; \alpha = 1)}{p_{\mathcal{R}}(y, z | a = 1; \alpha = 1)} \right). \quad (\text{A.7})$$

Statement (a) then directly follows from this equation. \square

A.4 Omitted Proofs from Subsection 5.1

We first derive the *IC* under the objective model \mathcal{R}^* . The probabilities of high output after high and low effort, respectively, are given by

$$p(y_H | a = 1) = \beta_4 + [\beta_2 + \beta_{02} + (\beta_1 + \beta_{01})\beta_{12}]\beta_{24} + [\beta_3 + (\beta_1 + \beta_{01})\beta_{13} + (\beta_2 + \beta_{02} + (\beta_1 + \beta_{01})\beta_{12})\beta_{23}]\beta_{34}, \quad (\text{A.8})$$

$$p(y_H | a = 0) = \beta_4 + [\beta_2 + \beta_1\beta_{12}]\beta_{24} + [\beta_3 + \beta_1\beta_{13} + (\beta_2 + \beta_1\beta_{12})\beta_{23}]\beta_{34}, \quad (\text{A.9})$$

so that the effect of effort on the probability of high output equals

$$p(y_H | a = 1) - p(y_H | a = 0) = (\beta_{02} + \beta_{01}\beta_{12})(\beta_{24} + \beta_{23}\beta_{34}) + \beta_{01}\beta_{13}\beta_{34}. \quad (\text{A.10})$$

Next, we derive the *IC* under the subjective model \mathcal{R} when the equilibrium action is $\alpha \in [0, 1]$.

We calculate

$$p(x_1 = 1 | x_2 = 1) = \frac{\alpha(\beta_1 + \beta_{01})(\beta_2 + \beta_{02} + \beta_{12}) + (1 - \alpha)\beta_1(\beta_2 + \beta_{12})}{\beta_2 + \beta_1\beta_{12} + \alpha(\beta_{02} + \beta_{01}\beta_{12})}, \quad (\text{A.11})$$

$$p(x_1 = 1 | x_2 = 0) = \frac{\alpha(\beta_1 + \beta_{01})(1 - \beta_2 - \beta_{02} - \beta_{12}) + (1 - \alpha)\beta_1(1 - \beta_2 - \beta_{12})}{1 - \beta_2 - \beta_1\beta_{12} - \alpha(\beta_{02} + \beta_{01}\beta_{12})}, \quad (\text{A.12})$$

and

$$p(x_3 = 1 | x_2 = 1) = \beta_3 + p(x_1 = 1 | x_2 = 1)\beta_{13} + \beta_{23}, \quad (\text{A.13})$$

$$p(x_3 = 1 | x_2 = 0) = \beta_3 + p(x_1 = 1 | x_2 = 0)\beta_{13}. \quad (\text{A.14})$$

The agent's belief about the probability of high output after $x_2 = 1$ and $x_2 = 0$, respectively, is therefore given by

$$p(y_H | x_2 = 1) = \beta_4 + \beta_{24} + [\beta_3 + p(x_1 = 1 | x_2 = 1)\beta_{13} + \beta_{23}]\beta_{34}, \quad (\text{A.15})$$

$$p(y_H | x_2 = 0) = \beta_4 + [\beta_3 + p(x_1 = 1 | x_2 = 0)\beta_{13}]\beta_{34}. \quad (\text{A.16})$$

The agent correctly anticipates $p(x_2 | a)$. Hence, her belief about the effect of effort on the probability of high output under \mathcal{R} equals

$$p_{\mathcal{R}}(y_H | a = 1; \alpha) - p_{\mathcal{R}}(y_H | a = 0; \alpha) = (\beta_{02} + \beta_{01}\beta_{12})(\beta_{24} + \beta_{23}\beta_{34}) + (\beta_{02} + \beta_{01}\beta_{12})\beta_{13}\beta_{34} \times [p(x_1 = 1 | x_2 = 1) - p(x_1 = 1 | x_2 = 0)]. \quad (\text{A.17})$$

Recall that $\beta_{13} < 0$. By comparing (A.10) and (A.17) we get that at $\alpha = 1$ the misspecification in \mathcal{R} relaxes the IC if and only if

$$\beta_{01} > \frac{\beta_{12}(\beta_1 + \beta_{01})(1 - \beta_1 - \beta_{01})(\beta_{02} + \beta_{01}\beta_{12})}{(1 - \beta_2 - \beta_{02} - \beta_{12}(\beta_1 + \beta_{01}))(\beta_2 + \beta_{02} + \beta_{12}(\beta_1 + \beta_{01}))}, \quad (\text{A.18})$$

which implies the statement in the main text.

Proof of Proposition 5. We prove the statements in (a). Since $\beta_1 \in \{0, 1\}$, we can rewrite the probability model without variable 1. The corresponding objective model $\tilde{\mathcal{R}}^*$ equals \mathcal{R}^* in Figure 4a without node 1. We now apply Propositions 7 and 8. In model $\tilde{\mathcal{R}}^*$, node 3 is not on a fundamental active path. Hence, the agent with subjective model \mathcal{R} is behaviorally rational, which yields the results. We prove the statements in (b). The first statement is shown in the text. The second statement follows from Corollary 3. Note that, in all models of Figure 4b, the set of nodes on fundamental active paths is identical. \square

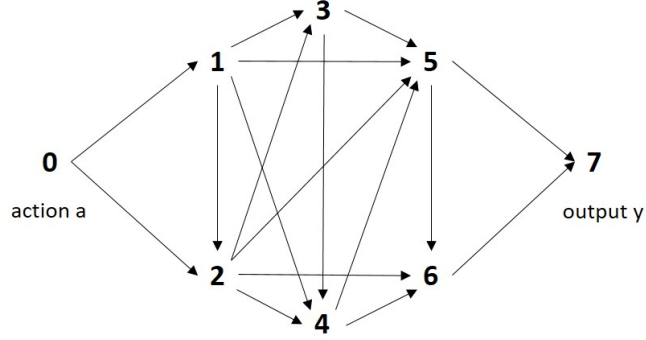
A.5 Omitted Proofs from Subsection 5.2

To prove Proposition 7, we first state and prove Proposition 8 below.

Proposition 8 (Fundamental Links). *Let \mathcal{R}^* be a perfect DAG and consider two adjacent nodes $i, j \in N^*$. The link iR^*j is fundamental if and only if at least one of the following conditions is satisfied:*

- (a) we have $d(0, i) = d(0, j) - 1$;
- (b) there exists a node $k \in N^*$ such that kEi and $k \notin R^*(j)$.

This result shows that nodes that are connected by fundamental links in perfect DAGs exhibit characteristics that are easy to identify. It is not always simple to spot the nodes that are not in $H^*(\mathcal{R}^*)$. In this case, Proposition 8 is helpful. Consider, for example, the perfect DAG \mathcal{R}^* in Figure 5. Condition (a) from Proposition 8 implies that all links which connect nodes of different distances to the action node are fundamental. The remaining links are $1R^*2$, $3R^*4$, $3R^*5$, $4R^*5$, $4R^*6$, and $5R^*6$. Condition (b) from Proposition 8 then implies that $4R^*6$ and $5R^*6$ are fundamental links, while the remaining links are non-fundamental. We therefore get $H^*(\mathcal{R}^*) = N^* \setminus \{3\}$.

Figure 5: Example DAG \mathcal{R}^* .

In order to prove Proposition 8, we show several intermediate results. We first note that in a perfect DAG \mathcal{R}^* the link iR^*j is fundamental if the nodes i and j differ in their distance to the action node 0.

Lemma 1. *Let $i, j \in N^*$ be adjacent nodes in \mathcal{R}^* . If $d(0, i) = d(0, j) - 1$, then iEj .*

Proof. First, suppose $d(0, i) = 0$ so that $i = 0$. Since node 0 is ancestral, we must have iGj in every DAG $\mathcal{G} \in \mathcal{E}$. Next, suppose $d(0, i) = d > 0$. Since \mathcal{R}^* is perfect and node 0 is ancestral, there exists an active path of length d from node 0 to node i . Denote by k the direct ancestor of i on this path. There cannot exist a link between k and j , otherwise we would have $d(0, i) = d(0, k)$, a contradiction. Thus, we must have iGk in every DAG $\mathcal{G} \in \mathcal{E}$, otherwise we would have a v -collider at node i . \square

Lemma 2. *Let $i, j \in N^*$ and iR^*j . If there exists a node $k \in N^*$ such that kEi and $k \notin R^*(j)$, then iEj .*

Proof. If there is a fundamental link from node k to node i , then iR^*j implies that we cannot have jR^*k . Otherwise, we would have a directed cycle. Node j and node k are therefore not adjacent. Hence, if jGi in some DAG $\mathcal{G} \in \mathcal{E}$, there would be a v -collider at i , a contradiction. \square

The “if”-statement of Proposition 8 follows directly from Lemma 1 and Lemma 2. For the “only if”-statement we need two more results. The first one provides a condition under which a link is not fundamental.

Lemma 3. *Let $i, j \in N^* \setminus \{0\}$ and iR^*j . If $R^*(i) \subset R^*(j)$, then the link between i and j is not fundamental.*

Proof. Consider the DAG $\mathcal{G} = (G, N^*)$ that is identical to \mathcal{R}^* except that it reverses the link between i and j . The assumption $R^*(i) \subset R^*(j)$ rules out that there are v -colliders in \mathcal{G} . Assume that there is a cycle in \mathcal{G} . Since \mathcal{R}^* is acyclic, the cycle must contain jGi . Further, there must exist a node k and a link kGj which is part of the cycle. Since \mathcal{R}^* is perfect, we must have $k\tilde{R}^*i$. Assume first that we have kR^*i . Then jGi implies that kGi is not part of the cycle. Thus, there must exist an active path τ of some length d so that $\tau_0 = i$ and $\tau_d = k$. But then there is a cycle consisting of the link kGi and τ . This cycle also exists in \mathcal{R}^* , a contradiction. Next, assume that we have iR^*k . Since $i \neq 0$ and $R^*(i) \subset R^*(j)$, there exists a node l with lR^*i and lR^*j . Since \mathcal{R}^* is perfect, we also must have $l\tilde{R}^*k$. The same applies to all $l' \in R^*(i)$. Hence, starting from \mathcal{R}^* , we can reverse the links between i and j as well as between i and k and obtain a DAG $\mathcal{G}' \in \mathcal{E}$. \square

The second result needed for the proof of the “only if”-statement of Proposition 8 demonstrates that for each node i in a perfect DAG \mathcal{R}^* there exists a DAG $\mathcal{G} \in \mathcal{E}$ in which there is no non-fundamental link that points to i .

Lemma 4. *For all nodes $i \in N^*$ there exists a DAG $\mathcal{G} \in \mathcal{E}$ in which all non-fundamental links adjacent to node i point away from i .*

Proof. Let N_d be the set of nodes that have distance $d > 0$ to the action node 0. Denote by $N_d^{[\kappa]}$, $\kappa = 1, 2, \dots$, the maximal subset of nodes that (i) are at distance $d > 0$ from the action node 0, and (ii) are connected through non-fundamental links (i.e., for any two nodes $i, j \in N_d^{[\kappa]}$ there exists a path between i and j consisting of non-fundamental links). **Step 1.** We show that all nodes in a given set $N_d^{[\kappa]}$ have the same parents outside of $N_d^{[\kappa]}$. Consider two nodes $i, j \in N_d^{[\kappa]}$ that are connected through the non-fundamental link iR^*j . By definition, we have kEi for each $k \in R^*(i) \setminus N_d^{[\kappa]}$ for each $i \in N_d^{[\kappa]}$. Since \mathcal{R}^* is perfect, this implies that $R^*(j) \setminus N_d^{[\kappa]} \subset R^*(i) \setminus N_d^{[\kappa]}$. Since iR^*j is non-fundamental, we also must have $R^*(i) \setminus N_d^{[\kappa]} \subset R^*(j) \setminus N_d^{[\kappa]}$ so that $R^*(i) \setminus N_d^{[\kappa]} = R^*(j) \setminus N_d^{[\kappa]}$. The result follows from the fact that, by assumption, all nodes in $N_d^{[\kappa]}$ are connected through non-fundamental links. **Step 2.** Consider two links $i \in N_d^{[\kappa]}$ and $i' \in N_d^{[\kappa']}$ with $\kappa \neq \kappa'$ that are adjacent. Assume w.l.o.g. that iR^*i' . By definition, iR^*i' is a fundamental link. Step 1 then implies that iEj' for all $j' \in N_d^{[\kappa']}$. Thus, there cannot exist nodes $j \in N_d^{[\kappa]}$ and $j' \in N_d^{[\kappa']}$ so that $j'R^*j$. Otherwise, we would have $j'Ej$ and $j'Ei$ for all $i \in N_d^{[\kappa]}$, a contradiction. Thus, there cannot exist nodes $i, j \in N_d^{[\kappa]}$ and $i', j' \in N_d^{[\kappa']}$ such that iR^*i' and $j'R^*j$. **Step 3.** Note that, since \mathcal{R}^* is perfect, by Lemma 1 all links between N_d and N_{d+1} point away from the nodes in N_d . **Step 4.** We now can prove Lemma 4. Take any node $i \in N^*$ and assume w.l.o.g. that $i \in N_d^{[\kappa]}$. Consider the DAG $\mathcal{G}^{[\kappa]} = (N_d^{[\kappa]}, G^{[\kappa]})$ where $G^{[\kappa]}$ is identical to \mathcal{R}^* restricted on $N_d^{[\kappa]}$. Since \mathcal{R}^* is perfect, $\mathcal{G}^{[\kappa]}$ also must be perfect. Corollary 1 from Spiegler (2019) implies that there exists a DAG $\mathcal{Q}^{[\kappa]}$ in which node i is ancestral and that is equivalent

to $\mathcal{G}^{[k]}$. Choose such a $Q^{[k]}$ and replace $\mathcal{G}^{[k]}$ in the original DAG \mathcal{R}^* by $Q^{[k]}$. Call the resulting DAG Q^* . Step 1 implies that there are no ν -colliders in Q^* , and Step 2 and 3 imply that there are no cycles in Q^* , which proves the result. \square

Proof of Proposition 8. The “if”-statement follows from Lemma 1 and Lemma 2. We prove the “only if”-statement. Consider any two adjacent nodes $i, j \in N^*$ with iR^*j and $d(0, i) = d(0, j)$. Suppose that for any node $k \in R^*(i)$ with a fundamental link kR^*i we also have $k \in R^*(j)$. By Lemma 4, we can find a DAG $\mathcal{G} \in \mathcal{E}$ in which all non-fundamental links are turned away from node i . In this DAG, we have $G(i) \subset G(j)$. From Lemma 3 it then follows that the link iR^*j is not fundamental. This completes the proof. \square

Before we can prove Proposition 7, we need two more results. We will use the following definitions. Recall that a path τ of length d is directed if for any $h \in \{1, \dots, d\}$ we have $\tau_{h-1}R\tau_h$ on this path. For any DAG, the topological ordering is a sequence of nodes such that every link is directed from an earlier to a later node in the sequence.

Lemma 5. *Let $M \subset N^* \setminus H^*(\mathcal{R}^*)$ be a set of nodes connected through non-fundamental links. Suppose there are two nodes $i, j \in H^*(\mathcal{R}^*)$ with non-fundamental links to nodes in M . Then i and j are adjacent.*

Proof. As in the proof of Lemma 4, let N_d be the set of nodes that have distance $d > 0$ to the action node 0. Let $E(i)$ be the set of nodes k with kEi . By Lemma 1, there is a $d > 0$ so that $i, j \in N_d$ and $M \subset N_d$. By Lemma 2, we must have $E(i) = E(j)$ since these nodes are connected through non-fundamental links. Choose any node $k \in N_{d-1}$ with $k \in H^*(\mathcal{R}^*)$ and kR^*i . By Lemma 2, we also have kR^*j . We can now choose two fundamental active paths $\tau^{[i]}, \tau^{[j]}$ from node 0 to node n so that (i) $k \in \tau^{[i]}$ and $k \in \tau^{[j]}$, (ii) $i \in \tau^{[i]}$ and $j \in \tau^{[j]}$, (iii) all nodes on $\tau^{[i]}$ and $\tau^{[j]}$ before k are identical, and (iv) there is not any node on $\tau^{[i]}$ ($\tau^{[j]}$) between k and i (k and j). Since $i, j \in H^*(\mathcal{R}^*)$ this is possible. Now define by $m_1^{[i]}$ ($m_1^{[j]}$) the last node on $\tau^{[i]}$ ($\tau^{[j]}$) before node n ; by $m_2^{[i]}$ ($m_2^{[j]}$) the penultimate node on $\tau^{[i]}$ ($\tau^{[j]}$) before node n , and so forth. Since \mathcal{R}^* is perfect, $m_1^{[i]}$ and $m_1^{[j]}$ must be adjacent. Since $m_1^{[i]}$ and $m_1^{[j]}$ are adjacent and \mathcal{R}^* is perfect, $m_2^{[i]}$ and $m_2^{[j]}$ must be adjacent, and so forth. If nodes i and j are both the t 'th node from n in $\tau^{[i]}$ ($\tau^{[j]}$), we are done. Assume that this is not the case, and that w.l.o.g. node i is the t 'th node from n while node j is the t' 'th node from n , with $t' > t$. Then i is adjacent to $m_t^{[j]}$, and also to all nodes on $\tau^{[j]}$ between $m_t^{[j]}$ and j (including j) through non-fundamental links, otherwise there would be a contradiction to $E(i) = E(j)$. \square

The next result is crucial for the proof of Proposition 7. It shows that all nodes that are not on a fundamental active path between action and output can be made “unimportant”, in the

sense that we can find a DAG in \mathcal{E} in which any link between a node in $H^*(\mathcal{R}^*)$ and a node in $N^* \setminus H^*(\mathcal{R}^*)$ points towards the node in $N^* \setminus H^*(\mathcal{R}^*)$.

Lemma 6. *There exists a DAG $\mathcal{G}^* \in \mathcal{E}$ such that in \mathcal{G}^* all links with one end in $H^*(\mathcal{R}^*)$ and the other in $N^* \setminus H^*(\mathcal{R}^*)$ point from $H^*(\mathcal{R}^*)$ to $N^* \setminus H^*(\mathcal{R}^*)$.*

Proof. The proof proceeds by steps. **Step 1.** Consider any maximal set $M \subset N^* \setminus H^*(\mathcal{R}^*)$ of nodes connected through non-fundamental links and let $M^+ \subset H^*(\mathcal{R}^*)$ be the set of nodes that have non-fundamental links to nodes in M . By Lemma 1, there is a $d > 0$ so that $M, M^+ \subset N_d$. Denote by M^{++} the set of nodes in $N_d \cap H^*(\mathcal{R}^*)$ with fundamental links into M . Since the nodes in M are connected through non-fundamental links, there is a fundamental link from any node $i \in M^{++}$ to any node in M . Thus, any node in M^{++} must also be adjacent to any node in M^+ , so $M^+ \cup M^{++}$ is a clique. **Step 2.** Consider the DAG $\bar{\mathcal{G}} = (N, \bar{G})$, where $N = M \cup M^+ \cup M^{++}$ and \bar{G} is identical to \mathcal{R}^* restricted on N . By construction, this DAG is perfect. Hence, Corollary 1 from Spiegler (2019) implies that there exists a DAG $\bar{\mathcal{G}}^+$ in which the clique $M^+ \cup M^{++}$ is ancestral and that is equivalent to $\bar{\mathcal{G}}$. We choose such a $\bar{\mathcal{G}}^+$ with the property that the ordering of the nodes in $M^+ \cup M^{++}$ is the same as in $\bar{\mathcal{G}}$ (this is possible since $M^+ \cup M^{++}$ is a clique, and all links between nodes $M^+ \cup M^{++}$ and nodes in M point towards the latter one). Consider now the DAG \mathcal{G} that is identical to \mathcal{R}^* except that $\bar{\mathcal{G}}$ is replaced by $\bar{\mathcal{G}}^+$. We show that there are no cycles or ν -colliders in \mathcal{G} so that it is equivalent to \mathcal{R}^* . Consider any node $i \in N_{d-1} \cup N_d$ that is outside $M \cup M^+ \cup M^{++}$ and that has a fundamental link into a node in M . Since the nodes in M are connected through non-fundamental links, node i has a fundamental link into every node in M (otherwise, i would belong to M , a contradiction). This rules out ν -colliders. Any link between a node in N_d and a node in N_{d+1} points into the latter one. Hence, by construction, there cannot be cycles or ν -colliders in \mathcal{G} . We obtain \mathcal{G}^* by performing the same changes for any maximal set $M \subset N^* \setminus H^*(\mathcal{R}^*)$ of nodes connected by non-fundamental links in \mathcal{R}^* . \square

Proof of Proposition 7. First, we show the “if”-statement. Assume that the agent’s subjective model \mathcal{R} contains all the nodes in $H^*(\mathcal{R}^*)$. Consider the DAG $\mathcal{G}^* \in \mathcal{E}$ in which all links with one end in $H^*(\mathcal{R}^*)$ and the other in $N^* \setminus H^*(\mathcal{R}^*)$ point from $H^*(\mathcal{R}^*)$ to $N^* \setminus H^*(\mathcal{R}^*)$. By Lemma 6, this DAG exists. From Proposition 6 it follows that $p_{\mathcal{G}^*}(x_{H^*(\mathcal{R}^*)}) = p(x_{H^*(\mathcal{R}^*)})$ for all distributions $p(x) \in \Delta(X)$. Consider the subgraph $\mathcal{G} = (G, N)$ where G equals \mathcal{G}^* restricted on N . Since none of the nodes in $N \setminus H^*(\mathcal{R}^*)$ impacts on any node in $H^*(\mathcal{R}^*)$, we have $p_{\mathcal{G}}(x_{H^*(\mathcal{R}^*)}) = p_{\mathcal{G}^*}(x_{H^*(\mathcal{R}^*)})$ for all $p(x) \in \Delta(X)$. By construction, the DAGs \mathcal{R} and \mathcal{G} are equivalent so that we have $p_{\mathcal{R}}(x_{H^*(\mathcal{R}^*)}) = p_{\mathcal{G}}(x_{H^*(\mathcal{R}^*)}) = p_{\mathcal{G}^*}(x_{H^*(\mathcal{R}^*)}) = p(x_{H^*(\mathcal{R}^*)})$ for all distributions $p(x) \in \Delta(X)$, which proves the “if”-statement. Next, we show the “only if”-statement. Assume that there is one node $i \in H^*(\mathcal{R}^*)$ that is not in the agent’s subjective model. This node is on a fundamental active path τ between the action node 0 and the output node n . We then can find a probability

distribution $p(x) \in \Delta(X)$ so that $p_{\mathcal{R}}(x_n | x_0) \neq p(x_n | x_0)$. Let k be the k 'th node in τ . Consider a probability distribution with the following properties: $p(x_j | x_{R^*(j)}) = p(x_j)$ for all nodes $j \notin \tau$ that are between the nodes 0 and n , and $p(x_k | x_{R^*(k)}) = p(x_k | x_{k-1})$. Clearly, such a distribution can have the desired property. \square

Proof of Corollary 3. Denote $H^*(\mathcal{R}_1) = H^*(\mathcal{R}_2) = H$. By Proposition 7, there exists a DAG $\mathcal{R}_1^{[1]}$ that is equivalent to \mathcal{R}_1 and in which all links between any node $i \in H$ and any node $j \in N_1 \setminus H$ is turned away from i . Thus, we have

$$p_{\mathcal{R}_1}(x_H) = \sum_{x_{N_1 \setminus H} \in X_{N_1 \setminus H}} p_{\mathcal{R}_1}(x_{N_1}) = \sum_{x_{N_1 \setminus H} \in X_{N_1 \setminus H}} p_{\mathcal{R}_1^{[1]}}(x_{N_1}) = p_{\mathcal{R}_1^{[1]}}(x_H). \quad (\text{A.19})$$

Note that for all $i \in H$ we have that $R_1^{[1]}(i) \subset H$. Consider the restriction of $R_1^{[1]}$ on H , $R_1^{[H]}$. We then have

$$p_{\mathcal{R}_1^{[1]}}(x_H) = \prod_{i \in H} p(x_i | x_{R_1^{[1]}(i)}) = \prod_{i \in H} p(x_i | x_{R_1^{[H]}(i)}) = p_{\mathcal{R}_1^{[H]}}(x_H). \quad (\text{A.20})$$

Define $\mathcal{R}_2^{[1]}$ and $\mathcal{R}_2^{[H]}$ just like $\mathcal{R}_1^{[1]}$ and $\mathcal{R}_1^{[H]}$. By assumption, the link $iR_1^{[H]}j$ is in $R_1^{[H]}$ if and only if we have $iR_2^{[H]}j$ or $jR_2^{[H]}i$. Thus, $\mathcal{R}_1^{[H]}$ and $\mathcal{R}_2^{[H]}$ have the same skeleton. Since \mathcal{R}_1 and \mathcal{R}_2 are perfect, so are $\mathcal{R}_1^{[H]}$ and $\mathcal{R}_2^{[H]}$. Hence $\mathcal{R}_1^{[H]}$ and $\mathcal{R}_2^{[H]}$ are equivalent, so that

$$p_{\mathcal{R}_1^{[H]}}(x_H) = p_{\mathcal{R}_2^{[H]}}(x_H). \quad (\text{A.21})$$

From the equations (A.19) to (A.21), we get $p_{\mathcal{R}_1}(x_H) = p_{\mathcal{R}_2}(x_H)$, which implies the result. \square