

# **Essays in Applied Microeconomics**

Inaugural-Dissertation  
zur Erlangung des Grades eines Doktors  
der Wirtschafts- und Gesellschaftswissenschaften  
durch die  
Rechts- und Staatswissenschaftliche Fakultät  
der Rheinischen Friedrich-Wilhelms-Universität  
Bonn

vorgelegt von  
Venuga Yokeeswaran  
aus Jaffna

Bonn 2015

Dekan:	Prof. Dr. Rainer Hüttemann
Erstreferent:	Prof. Dr. Dezsö Szalay
Zweitreferent:	Prof. Dr. Matthias Kräkel

Tag der mündlichen Prüfung: 21.09.2015

# Acknowledgments

This PhD thesis is the result of a journey I wouldn't have completed without the contribution and support of many people around me. I am deeply indebted to each and every one of them.

First, I would like to thank my first supervisor, Dezső Szalay, for providing his time and beneficial comments throughout my whole PhD and for giving me the freedom to work on projects from diverse areas of economics. I am also very grateful to my second supervisor, Matthias Kräkel, for his precious comments and his kind support.

I would like to express my deepest gratitude to my three co-authors Dezső Szalay, Mark Le Quement, and Renaud Coulomb. All three collaborations were inspiring, motivating and supportive. It was a pleasure to work with them.

During my research stay at the London School of Economics (LSE) as part of the European Doctoral Program (EDP) I received valuable support from Ronny Razin for which I am sincerely thankful.

Steffen Altmann, Mark Le Quement, Felix Pasker, and Matthias Wibrals deserve special mention for proofreading and very useful advice on draft versions of this dissertation and their constant support and encouragement.

I am much obliged to all my friends for the valuable, pleasant, and inspiring time together and for their moral support when it was most needed.

Financial support from the German Research Foundation (DFG) and the Bonn Graduate School of Economics (BGSE) is gratefully acknowledged.

Last but not least, I am deeply grateful to my beloved family and my dearly-loved partner for enabling and supporting my studies and for their continuous encouragement and support in every matter during the challenging phases.



# Contents

<b>Introduction</b>	<b>1</b>
<b>1 Managerial Incentive Problems and Return Distributions</b>	<b>5</b>
1.1 Introduction . . . . .	5
1.2 The Model . . . . .	10
1.3 The Principal's Problem . . . . .	12
1.4 The Problem of Pure Moral Hazard . . . . .	14
1.4.1 Optimal Contracts . . . . .	15
1.4.2 Covariance of Contracts and Moments of the Profit Distribution . . . . .	17
1.5 The Case of Combined Adverse Selection and Moral Hazard . . . . .	18
1.5.1 Optimal Contracts . . . . .	21
1.5.2 Covariance of Contracts and Moments . . . . .	23
1.6 Attenuation . . . . .	23
1.7 Conclusions . . . . .	25
<b>2 Subgroup Deliberation and Voting</b>	<b>29</b>
2.1 Introduction . . . . .	29
2.2 The Model . . . . .	32
2.2.1 Setup . . . . .	32
2.2.2 Communication Protocols and Equilibria . . . . .	34
2.3 Positive Analysis . . . . .	40
2.4 Normative Analysis . . . . .	42
2.5 Conclusion . . . . .	47

<b>3</b>	<b>Carbon Taxation under Asymmetric Information over Fossil-fuel Reserves</b>	<b>51</b>
3.1	Introduction . . . . .	51
3.2	The Model . . . . .	56
3.3	The Case of Symmetric Information . . . . .	60
3.4	The Case of Asymmetric Information . . . . .	64
3.4.1	Separating Equilibria . . . . .	66
3.4.2	Pooling Equilibria . . . . .	70
3.4.3	Equilibrium Selection . . . . .	76
3.5	Welfare Analysis . . . . .	78
3.6	Conclusion . . . . .	80
	<b>Appendix</b>	<b>83</b>
<b>A</b>	<b>Managerial Incentive Problems and Return Distributions</b>	<b>83</b>
A.1	The Problem of Pure Moral Hazard . . . . .	83
A.2	The Case of Combined Adverse Selection and Moral Hazard . . . . .	87
<b>B</b>	<b>Subgroup Deliberation and Voting</b>	<b>97</b>
B.1	The Model . . . . .	97
B.2	Positive Analysis . . . . .	100
B.3	Normative Analysis . . . . .	107
<b>C</b>	<b>Carbon Taxation under Asymmetric Information over Fossil-fuel Reserves</b>	<b>117</b>
C.1	The Case of Symmetric Information . . . . .	117
C.2	The Case of Asymmetric Information . . . . .	123
C.3	Welfare Analysis . . . . .	134
	<b>Bibliography</b>	<b>139</b>

# Introduction

This thesis consists of three independent chapters, each covering a significant research field in applied microeconomics. A central part of economic studies is the problem of providing right incentives; e.g. when delegating tasks to employees the employer should make sure that the assignments are executed in his interest. The issue of not having the same goals arise when contracting partners have conflicting objectives. If, in addition, the employee's characteristics are not accurately known to the employer, if the tasks are to be taken in a risky environment, and if the employee is risk-averse, the problem of finding right incentives becomes more relevant, but, also more complex. The literature on incentive theory, or contract theory, deals with these questions and problems. The first chapter of this thesis is a contribution to this literature and to related empirical studies. A rather general principal-agent model is used to address empirical findings concerning risk-incentive trade-offs.<sup>1</sup> Another substantial field in economics is political economy. The second and third chapters of this thesis are contributions to the literature on the associated branches, political economy of collective decision making and political economy of climate change, respectively. Crucial questions in the first subfield are how to best aggregate private information and how collective decision making performs. In America, e.g., defendants are judged by a jury. Members of a jury might have diverging preferences and do have their own perception of the defendant's guiltiness. The questions of whether it is possible to extract or to determine each juror's true assessment and whether the jury reaches a verdict that does the defendant justice, is socially very important. The second chapter of this thesis compares jury verdicts under three different protocols for aggregating information and ranks them with respect to welfare.<sup>2</sup> Finally, the third chapter deals with environmental economics. Global warming and climate change are amongst the

---

<sup>1</sup>This chapter is based on the paper "Managerial Incentive Problems and Return Distribution" which is joint work with Dezső Szalay.

<sup>2</sup>This chapter is based on the paper "Subgroup Deliberation and Voting" which is joint work with Mark Le Quement and published in *Social Choice and Welfare* (Le Quement & Yokeswaran 2015).

major threats mankind will have to face in the future. It is therefore essential to study optimal regulation of carbon emissions and fossil-fuel consumption. The last chapter of this thesis tackles this problem. It analyzes optimal Pigovian taxation in a “delayed regulation” model with asymmetric information on fossil-fuel reserves. Furthermore, it makes predictions about the effects of the regulation on relevant parties.<sup>3</sup>

These three chapters are now described in more detail individually.

**Chapter 1** Contracting partners commonly have diverging interests. This makes it relevant to study optimal incentive schemes by the use of principal-agent models and to characterize the relationship between risk and incentives. Standard principal-agent theory (Holmström & Milgrom 1987) determines a negative trade-off between risk and incentives. A principal, offering a performance pay to a risk-averse agent, trades off the benefit from higher effort to the loss from higher risk compensation. In a high risk environment the performance pay is therefore low. Various empirical studies testing this theory find conflicting results on whether to support this prediction or not. Some papers indeed find a negative relationship, others claim that the relationship is positive while there exist studies that do not find a significant relationship at all (see e.g. Aggarwal & Samwick 2002, Core & Guay 1999, Bushmann et al. 1996). Successful efforts have already been made to explain the positive relationship (see e.g. Prendergast 2002). The main idea of these papers is to incorporate aspects of managerial incentive problems which are neglected in the standard model such as e.g. the possibility of endogenous delegation of decisions, or of endogenous matching. Prendergast (2002), who looks into the first mentioned extension, assumes that managers are better informed on output, hence have greater value in a high risk environment, and therefore should be highly incentivised to manage in the principal’s interest. The positive relationship is then a natural outcome if the data and the tested industry are liable to these aspects. Our contribution to this literature is to incorporate the issue of the agent being able to choose endogenously the first two moments of the firm’s profit distribution. We study a general model and make comprehensive comparative statics predictions. We also explicitly emphasize the impact of endogenous risk on the theoretical predictions and the empirical evidence.

In our model the manager chooses the mean and the volatility of the firm’s profit distribution along an efficient frontier. The managers differ in two aspects, their cost of effort

---

<sup>3</sup>This chapter is based on the paper “Carbon Taxation under Asymmetric Information over Fossil-fuel Reserves” and is joint work with Renaud Coulomb.



and their risk aversion. If these characteristics are commonly known and associated, the relationship between volatility of profits and incentives is positive. Allowing for asymmetric information on these parameters the correlation stays positive, as long as the variation in the observed contracts is not too large. Consequently, our model also allows for negative correlation. In addition, we point out that empirical studies neglecting the endogeneity of risk –if the risk in the data is indeed endogenous– might falsely reject a significant relationship as this negligence biases estimates towards zero.

**Chapter 2** Early contributions to the literature on collective decision making compare different voting rules under private voting. Allowing for strategic voting the unanimity rule is known to aggregate private information poorly (see e.g. Feddersen & Pesendorfer 1998). Each juror acts as if his vote is pivotal when voting strategically. Being pivotal reveals additional information about the defendant’s guiltiness which might overwhelm the juror’s own assessment. To improve on information aggregation various papers add a communication stage prior to voting. Most studies restrict communication to simultaneous plenary deliberation and study the possibility of truthful deliberation and sincere voting. Coughlan (2000), however, shows that truthful deliberation does not constitute an equilibrium outcome if committee members have commonly known and substantially heterogeneous preferences. Jurors might not want to reveal their true assessment in order to manipulate different minded jurors. Austen-Smith & Feddersen (2006) show that uncertainty about these preferences can render full pooling of information compatible with heterogeneity, as long as the voting rule is not unanimity. Gerardi & Yariv (2007) generalize the communication and the voting stages by not specifying the communication and the voting protocols and show that all voting rules, but the unanimity rule, are equivalent. An outcome that can be implemented by one voting rule can also be implemented by any other voting rule, by agreeing on this outcome in the communication stage and subsequently voting unanimously in favor of it, as long as the voting rule is non-unanimous. Consequently, the jurors are never pivotal when voting and hence do not have an incentive to deviate. Our contribution to this literature is to take the poorly performing unanimity rule and to introduce a communication protocol different from plenary deliberation in order to compare the outcomes. Additionally, we emphasize that this new protocol can improve the outcome.

We consider a heterogeneous committee voting by unanimity rule. We treat three different protocols to aggregate information, private voting and voting preceded by either plenary or

subgroup deliberation. While the first deliberation protocol imposes public communication, the second one restricts communication to homogeneous subgroups. We find that both protocols allow to Pareto improve on outcomes achieved under private voting. In addition, we find that when focusing on simple equilibria under plenary deliberation, subgroup deliberation Pareto improves on outcomes achieved under plenary deliberation.

**Chapter 3** According to climate change experts dangerous climate change cannot be prevented without reducing fossil-fuel combustion (see e.g. IEA 2012). Regulators' responsibility is to intervene and find optimal measures to control fossil-fuel consumption and hereby regulate carbon emissions. Optimal carbon regulation in a world in which all information is publicly known is well researched in various settings based on the Hotelling model (see e.g. Ulph & Ulph 1994). However, it seems plausible that oil owners have private information on their reserves (IEA 2010); empirical studies such as Bentley (2002) and Laherrere (2013) find that reserves of non-renewable resources are commonly over-reported. The environmental literature has examined optimal taxation in different asymmetric information settings (see e.g. Jebjerg & Lando 1997, Osmundsen 1998). However, the interaction of carbon taxation and the revelation of private information on the size of the reserves has not been analyzed yet. Our contribution to this literature is to address this matter and to examine the optimal carbon taxation under asymmetric information on the size of the fossil-fuel reserves. Our results suggest that a threat of a future mandatory carbon taxation might even be an explanation for the evidence of the over-reporting mentioned above.

We examine a setting in which a delayed environmental regulation –in order to reflect slow international negotiations– is implemented by a Pigovian tax. The regulator aims to control the environmental damages caused by CO<sub>2</sub> emissions from fossil-fuel combustion while not neglecting the social welfare from its usage. Information on the size of the reserves –low or high– is the resource owner's private information. We find that the threat of carbon regulation creates incentives to exaggerate the size of the reserves. This behavior does not have to be detrimental to future environmental regulation. Both parties, the resource owner who wants to maximize his profits and the social planner who wants to control environmental pollution, can profit from asymmetric information.

## Chapter 1

# Managerial Incentive Problems and Return Distributions

### 1.1 Introduction

The separation of ownership and control (Berle & Means 1932) makes it vital to understand the optimal design of incentive schemes for managers. The theoretical literature on the subject is vast, leaving no hope to do justice to all contributions. A cornerstone of incentive theory is the Holmström & Milgrom (1987) continuous time model, where a manager's compensation takes the form of a linear compensation scheme; a fixed part plus a variable part that depends linearly on some accounting measure; Hellwig & Schmidt (2002) have provided discrete time approximations for this model. Applied work on contract theory usually starts from a static reduced form version of these models, assuming that a manager receives a compensation package that is linear in profits, and studies how the components of the manager's pay change with the underlying problem. One comparative statics prediction that is shared by the majority of these models is that the sensitivity of the manager's pay to the firm's profits should be the lower the more risky the firm's profits are. Efficient risk sharing between well diversified shareholders and the firm's managers would allocate all the risk to the shareholders, but such an arrangement would give the manager too little incentives to work. Hence, moral hazard induces an inefficiency that is the more costly the larger the underlying risk and so the optimal sensitivity of the manager's performance pay is reduced when the firm's profits become more volatile.

The empirical evidence as to whether the data support this comparative statics prediction is mixed. In the context of executive pay, Core & Guay (1999) and Oyer & Schaefer (2004) find a positive and significant relation between measures of business risk and performance sensitivity of pay; Aggarwal & Samwick (2002) and Lambert & Larcker (1987) find a negative and significant relation between risk and incentives. Quite some studies find results that are statistically not significant: Bushmann et al. (1996) and Ittner et al. (1997) study whether firms are more or less inclined to use individual performance evaluation rather than compensation based on financial performance measures when risk is higher and find a positive result when they take variance in stock returns as the measure of risk; they find a negative result when they take variance in accounting returns as their measure of risk; Ittner et al. (1997) find positive results for various measures of risk (volatility of accounting returns, stock returns and net earnings); Yermack (1995) finds that firms provide more incentives from stock options when accounting earnings contain larger amounts of noise.

We propose a new way to look at this evidence. We develop a theoretical model of performance pay where the manager is given incentives to be diligent in two respects. Firstly, the manager exerts effort which, all else equal, makes higher profits more likely. Secondly, the manager can also choose the firm's strategy, that is, he can select the riskiness of the firm's profits along an efficient frontier. We stick firmly to the applied perspective and assume that the manager faces a compensation package that is linear in profits<sup>1</sup>. The performance sensitivity of the manager's pay determines both his optimal effort choice and the optimal volatility of the firm's profits. The optimal contract is influenced by the manager's underlying characteristics. When these characteristics vary, the observed contract choices vary too and furthermore induce variation in the observed firm characteristics. Hence, our model makes predictions as to the covariation between observed contract choices and firm characteristics, that is, mean and variance of profits. Since we do not in general know whether the characteristics are known to the principals who design the contracts (in practice), we extend our results to allow for adverse selection with respect to the manager's characteristics on top of moral hazard with respect to the choices made by the manager. Under fairly general conditions, we obtain a (pairwise) positive covariance of performance-sensitivity of pay and mean returns

---

<sup>1</sup>This is a standard perspective taken in a sizeable branch of the literature. While restricting contracts to a particular functional form is clearly a restriction, doing so allows us to closely compare our results to those found in the applied literature that works from this hypothesis, which is precisely the aim of the present chapter. Thus, the restriction to linear contracts is imposed deliberately, not just for analytical convenience.

and volatility of profits.

If there is a grain of truth to our story, then our model sheds new light on the existing evidence. The hypothesis that risk and incentives should be inversely related is based on a model where risk is exogenous. In contrast, when risk is endogenous through choices made by the managers, then our theoretical model predicts a positive relation. Moreover, in empirical studies endogeneity would not only affect the sign but also the magnitude of the estimated relations, at least when the endogeneity is not entirely accounted for: the resulting correlations between risk as a regressor and the error terms biases the estimates towards zero, explaining why it is difficult to reject the hypothesis that there is no relation between risk and incentives at all.

Our story is closely related to Holmström & Milgrom (1994) and Demski & Dye (1999). Holmström & Milgrom (1994) develop a theory of the joint determination of various elements of contracts. While their theory explains the covariation of choices made by the principal, we wish to explain how choices made by principals (that is, contracts) covary with choices made by the managers (that is, expected level and riskiness of profits). A key element in our theory is an efficient frontier, which introduces a relation between equilibrium expected return and risk. Demski & Dye (1999) also build on the idea that a manager can make mean-variance trade-offs; however, they address quite different questions with their model.

Thus, the key idea is to allow for more margins of decision making that affect the contracting environment. This idea is also present in Hellwig (2009), Sung (2005), Araujo et al. (2007) and Garcia (2014). All these papers allow for endogenous risk choices, even though the precise trade-offs and who controls the choice of risk differs across the approaches. Hellwig (2009) points out that all moments, mean and risk, are jointly determined as solutions of one incentive problem and thus challenges the way we think of debt contracts as a solution to one incentive problem and equity contracts as a solution to another one. Sung (2005) studies a continuous time principal agent problem with moral hazard and adverse selection which allows for an endogenous choice of volatility by the principal. We use a static model but allow for various sources of heterogeneity among agent types, have all choices except for contracts made by the agent, and explore the comparative statics properties based on the association of random variables as Holmström & Milgrom (1994) do<sup>2</sup>. Araujo et al. (2007)

---

<sup>2</sup>Combining Sung's (2005) continuous time with our multidimensional approach is - as we believe - an interesting avenue for future research.

analyze a problem where the manager's effort choice raises means and reduces variance at the same time. In Garcia (2014), risk can again be seen as an additional contracting tool that the principal uses alongside with linear contracts to control the agent's effort choice.

Overall, we believe it is very natural to assume that all the moments of the return distribution are endogenous and find it reassuring that different variations on the same theme share similar results.<sup>3</sup> Many variations and their predictions for empirical work remain unexplored to date.<sup>4</sup>

Part of the empirical contracting literature discusses endogeneity of risk explicitly; see, e.g. Garen (1994) and, in the context of franchising, Lafontaine (1992), Lafontaine & Slade (1998), and Lafontaine & Slade (2007). One way to deal with the issue is to find measures of risk that are likely to be exogenous to the firm's choices. Garen (1994) follows this approach and uses R&D intensity as a proxy for the riskiness of a firm's industry. Using that proxy he finds a negative but statistically not significant relation between the pay-performance sensitivity and this proxy. Based on our theory, we propose an empirical approach that attacks this issue more directly, that is to regress all the choices made by the manager and the firm's owners on the characteristics of the underlying problem. This would allow to estimate the endogenous relation between risk and incentives.

While we are not aware of any study in the context of executive pay that addresses this issue, Akerberg & Botticini (2002) make a closely related point in the context of sharecropping, pointing out that some characteristics of their underlying contracting problem may be endogenous through tenant/landowner matching. In their context, the landowner decides on what crop to grow; if crops differ in their riskiness, then tenants who differ in their risk aversion feel attracted to different landowners. Similar to their work, we stress that endogeneity is an important issue. However, since the details of optimal choices in a contracting relationship are different from the details in the matching process, our way to address the endogeneity is quite different.

A number of theories can rationalize a positive relation between risk and incentives. The main value added to our exercise is not so much to provide yet another one explaining the

---

<sup>3</sup>It should also be stressed that incentive problems in practice may depend on the context, ranging from excessive risk taking to excessive conservatism. This chapter does not address excessive risk taking by managers, and is therefore clearly not the adequate framework to think about contracts for bank managers, where excessive risk taking is the main concern. For a further discussion, see the final section.

<sup>4</sup>In a recent paper, Weinschenk (2014) studies a model with endogenous project choice to challenge the Marshallian hypothesis that higher incentives lead to higher expected profits - a feature that our model has. He shows that this need not be the case in his model.

same thing but much more to paint a rich picture of the comparative statics predictions of a contracting model allowing for many margins along which managers make choices and for many dimensions of heterogeneity among managers that may or may not be private information within a unified framework. We are not aware of a similar attempt in the literature.

Prendergast (2002) was first to take up the mismatch between theory and empirical work. He argues that the standard theory neglects an endogenous delegation decision. Suppose there are two essential inputs in production, effort and information that is used to make decisions, and suppose that agents have better information than principals. The value of this improved information is the larger the more uncertain the environment. Consequently, the larger is business risk, the more likely are principals to delegate decision making to the agent. But to ensure that the agent acts in the principal's interest, the principal makes the agent's pay depend on his performance. Hence, the agent's pay is the more dependent on performance the higher is risk. Thus, essentially Prendergast (2002) argues that the existing theories and their empirical tests suffer from an omitted variable bias.

Raith (2003) argues that empirical tests of the principal agent model fail to distinguish variability in profits and measurement error in contracting. If this distinction is made, then a positive correlation of performance pay and business risk can be rationalized. In particular, he studies a model of oligopolistic competition, where a manager's role is to reduce his firm's costs of production. As in the traditional model, the dependence of the manager's pay on realized cost reductions is the smaller the larger is the measurement error for these same cost reductions. On the other hand, uncertainty about rivals' costs makes firms' profits stochastic. Although the power of managers' performance pay and the variability of firms' profits are not causally linked to each other, a change in a third factor, e.g. the degree of competition, increases both profit risk and the power of managers' incentive schemes. Thus, the agents' pay is more performance dependent when business risk is greater, but there is no causal link between the two effects.

More recently, Inderst & Müller (2010) point to the role of incentive pay-schemes when it comes to inducing exit by bad managers. Comparing severance pay with on the job payment schemes, they find that severance pay makes shirking too attractive for managers; on the other hand, risky pay-for-performance is only attractive to manager who think they are more likely to generate high returns. Moreover, performance pay may be steeper if the underlying

firm risk is higher.<sup>5</sup>

As stressed before, the main point of this chapter is not so much that the relation between risk and incentives is positive but that it is endogenous and shaped by many factors that may or may not be observed when contracts are written. We develop a framework which allows us to illustrate the implications of this insight for empirical work.

The remainder of this chapter is structured as follows. In section 1.2, we lay out the model. In section 1.3, we explain the principal's problem including its solution in the first-best situation. In section 1.4, we study the contracting problem with known characteristics, in section 1.5 we extend these results to the case of adverse selection with respect to the manager's characteristics. In section 1.6, we remind the reader of the attenuation problem in empirical studies that arises from endogeneity of regressors, in our case risk. Section 1.7 concludes. All proofs are gathered in the appendixes A.1 and A.2.

## 1.2 The Model

An owner of a firm hires a manager to produce output. Henceforth, we call the owner the principal (she) and the manager the agent (he). The distribution of profits,  $\pi$ , depends on the agent's management style, that is two choices the agent makes. In particular, the agent chooses the mean  $\mu$  and the variance  $\sigma^2$  of a Gaussian profit distribution, so  $\tilde{\pi} \sim N(\mu, \sigma^2)$ . The agent's choices are constrained by an efficient frontier  $\mu = \mu(e, \sigma)$ , where  $e$  is the agent's effort. The efficient frontier describes the maximum expected return the agent can reach for any given variance and effort choice. For a given effort, higher expected returns can only be reached at the cost of higher variance. By increasing his effort, the agent can expand the set of feasible profit distributions; the efficient frontier is increasing in  $e$  for any given volatility  $\sigma$ . We assume that  $\mu(e, \sigma)$  is jointly concave in  $e$  and  $\sigma$ . Finally, there is an upper bound on the volatility,  $\bar{\sigma}$ . Figure 1.1 depicts the efficient frontier.<sup>6</sup>

Effort is costly to the agent. The cost of effort is  $c \cdot e$ , where  $c$  is a positive parameter<sup>7</sup>. Contracts can only be written on profits; the agent's choices themselves - neither effort nor volatility - are not observable to the principal by the time payments are made. Moreover,

---

<sup>5</sup>A positive relation between risk and incentives can be rationalized in a number of other ways, e.g., through endogenous matching between principals and agents (Serfes 2005, Wright 2004) or by combining limited liability with risk aversion on the part of the agent (Budde & Kräkel 2011).

<sup>6</sup>For most of the chapter,  $\bar{\sigma}$  can be taken as  $\infty$ . Since the principal is risk neutral, we need a finite  $\bar{\sigma}$  to make the first-best allocation well defined.

<sup>7</sup>Making costs linear in effort is a normalization that is without loss of generality.



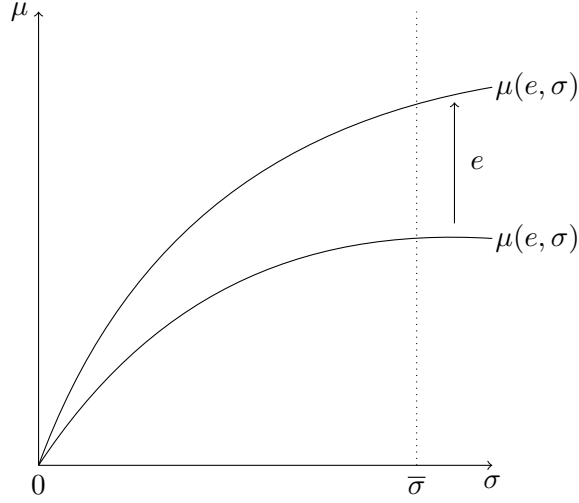


Figure 1.1: The efficient frontier

the principal is restricted to use linear contracts. So, the principal's wealth is equal to  $W_P = -\beta + (1 - \alpha)\pi$  and the agent's wealth is equal to  $W_A = \beta + \alpha\pi$ , where  $\beta$  is a base salary and  $\alpha$  the agent's share of profits. The principal is risk neutral while the agent is risk averse. His utility function displays constant absolute risk aversion. More precisely, we have

$$U_A(W_A, e) = -\exp(-a(W_A - ce)),$$

where  $a$  is the coefficient of absolute risk aversion. As is well known, the agent's expected utility can be expressed as  $\mathbb{E}[U_A(W_A, e)] = U_A(w_A - ce)$  where

$$w_A \equiv \beta + \alpha\mu(e, \sigma) - a\frac{\alpha^2}{2}\sigma^2 \tag{1.1}$$

is the agent's certainty equivalent level of wealth. Clearly, for any given effort choice, the agent will always choose a point on the efficient frontier, so  $\mu = \mu(e, \sigma)$ . The principal's expected utility is equal to his expected wealth, so

$$\mathbb{E}[W_P] = -\beta + (1 - \alpha)\mu(e, \sigma).$$

The agent's outside option gives rise to a certainty equivalent level of wealth of  $\omega$ . The agent knows his marginal cost of effort and his coefficient of risk aversion. These parameters are distributed with full support on the product set  $\mathbf{T} \equiv [\underline{a}, \bar{a}] \times [\underline{c}, \bar{c}]$  where  $\bar{a} > \underline{a} > 0$  and

$\bar{c} > \underline{c} > 0$ . We let  $\mathbf{t} \equiv (a, c)$  denote a type and let  $k(\mathbf{t})$  and  $K(\mathbf{t})$  denote the joint density and cdf of  $\mathbf{t}$ , respectively.

Apart from the efficient frontier - our key new element - these assumptions are standard in the literature (see e.g. Holmström & Milgrom 1994). We explore two variations of our model; in the first version, the agent's type is commonly known so that the only contractual friction is moral hazard arising from the unobservability of the agent's choices; in the second version, the principal only knows the distribution of the agent's type (and this is common knowledge), so there is adverse selection on top of moral hazard.

### 1.3 The Principal's Problem

We state the principal's problem for the most general case, where the agent has private information about his level of risk aversion and cost of effort. The case of symmetric information is then a special case of the general formulation.

Invoking the Revelation Principle, an optimal contract can be found restricting attention to a direct revelation game, where the agent is asked to announce his preference parameters  $\hat{\mathbf{t}}$ , and is given incentives to announce his type truthfully. For any given announced type,  $\hat{\mathbf{t}} \in \mathbf{T}$ , a contract specifies the quadruple  $\{\beta(\hat{\mathbf{t}}), \alpha(\hat{\mathbf{t}}), e(\hat{\mathbf{t}}), \sigma(\hat{\mathbf{t}})\}$ . Our problem is a combined problem of moral hazard and adverse selection. However, once the agent has announced a type,  $\beta(\hat{\mathbf{t}})$  and  $\alpha(\hat{\mathbf{t}})$  are given from his perspective. So, we can use (1.1) to compute the optimal choices of effort and standard deviation (from his perspective); let  $e(\alpha, \mathbf{t})$  and  $\sigma(\alpha, \mathbf{t})$  denote these choices. Since  $\mu(e, \sigma)$  is jointly concave in its arguments, incentive compatible choices are completely described by the pair of first-order conditions

$$\alpha(\hat{\mathbf{t}}) \mu_e(e(\alpha(\hat{\mathbf{t}}), \mathbf{t}), \sigma(\alpha(\hat{\mathbf{t}}), \mathbf{t})) = c \quad (1.2)$$

and either

$$\sigma(\alpha(\hat{\mathbf{t}}), \mathbf{t}) \leq \bar{\sigma} \text{ and } \mu_\sigma(e(\alpha(\hat{\mathbf{t}}), \mathbf{t}), \sigma(\alpha(\hat{\mathbf{t}}), \mathbf{t})) - a\alpha(\hat{\mathbf{t}}) \sigma(\alpha(\hat{\mathbf{t}}), \mathbf{t}) = 0, \quad (1.3)$$

or

$$\sigma(\alpha(\hat{\mathbf{t}}), \mathbf{t}) = \bar{\sigma} \text{ and } \mu_\sigma(e(\alpha(\hat{\mathbf{t}}), \mathbf{t}), \sigma(\alpha(\hat{\mathbf{t}}), \mathbf{t})) - a\alpha(\hat{\mathbf{t}}) \sigma(\alpha(\hat{\mathbf{t}}), \mathbf{t}) \geq 0. \quad (1.4)$$

Given strict concavity of the function  $\mu(\cdot, \cdot)$  in its arguments, this system of equations has a unique solution. Taking these choices into account, the principal's problem is reduced to a problem of pure adverse selection. The principal's problem is to

$$\max_{\beta(\cdot), \alpha(\cdot), \tilde{\mathbf{T}}(\omega)} \int_{\tilde{\mathbf{T}}(\omega)} (-\beta(\mathbf{t}) + (1 - \alpha(\mathbf{t})) \mu(e(\alpha(\mathbf{t}), \mathbf{t}), \sigma(\alpha(\mathbf{t}), \mathbf{t}))) k(\mathbf{t}) d\mathbf{t} \quad (1.5)$$

*s.t.*

$$w_A(\beta(\mathbf{t}), \alpha(\mathbf{t}), e(\alpha(\mathbf{t}), \mathbf{t}), \sigma(\alpha(\mathbf{t}), \mathbf{t})) - ce(\alpha(\mathbf{t}), \mathbf{t}) \quad (1.6)$$

$$\geq w_A(\beta(\hat{\mathbf{t}}), \alpha(\hat{\mathbf{t}}), e(\alpha(\hat{\mathbf{t}}), \mathbf{t}), \sigma(\alpha(\hat{\mathbf{t}}), \mathbf{t})) - ce(\alpha(\hat{\mathbf{t}}), \mathbf{t}) \text{ for all } \mathbf{t}, \hat{\mathbf{t}} \in \tilde{\mathbf{T}}(\omega)$$

$$w_A(\beta(\mathbf{t}), \alpha(\mathbf{t}), e(\alpha(\mathbf{t}), \mathbf{t}), \sigma(\alpha(\mathbf{t}), \mathbf{t})) - ce(\alpha(\mathbf{t}), \mathbf{t}) \geq \omega \text{ for all } \mathbf{t} \in \tilde{\mathbf{T}}(\omega) \quad (1.7)$$

$$\max_{\hat{\mathbf{t}} \in \tilde{\mathbf{T}}(\omega)} w_A(\beta(\hat{\mathbf{t}}), \alpha(\hat{\mathbf{t}}), e(\alpha(\hat{\mathbf{t}}), \mathbf{t}), \sigma(\alpha(\hat{\mathbf{t}}), \mathbf{t})) - ce(\alpha(\hat{\mathbf{t}}), \mathbf{t}) \leq \omega \text{ for all } \mathbf{t} \in \mathbf{T} \setminus \tilde{\mathbf{T}}(\omega) \quad (1.8)$$

In this problem, constraint (1.6) is the incentive constraint that guarantees truth-telling. (1.7) ensures that agents with characteristics  $\mathbf{t} \in \tilde{\mathbf{T}}(\omega)$  are willing to participate, (1.8) ensures that agents with other characteristics do not participate; the principal chooses the set  $\tilde{\mathbf{T}}(\omega)$ , i.e., whom to attract and whom to exclude. Note again that the moral hazard part of our problem has been subsumed into the hidden information part of the problem by requiring that  $e(\alpha(\mathbf{t}), \mathbf{t})$  and  $\sigma(\alpha(\mathbf{t}), \mathbf{t})$  satisfy the conditions (1.2) and either (1.3) or (1.4).<sup>8</sup> Thus, the problem of pure moral hazard corresponds to the problem above when we drop constraint (1.6). Moreover, in this case, the problem can always be solved pointwise for each  $\mathbf{t}$ .

The choice of the set  $\tilde{\mathbf{T}}(\omega)$  is only interesting in case the characteristics are privately known to the agent; this is due to the absence of wealth effects in the principal's and the agent's utility function. For this reason we study the problem with known characteristics under conditions that ensure that full participation is optimal; formally, we set  $\omega = 0$  for the first part of the analysis in section 1.4, which ensures that  $\tilde{\mathbf{T}}(\omega) = \mathbf{T}$ . In contrast, the optimal allocation in the case of privately known characteristics, analyzed in section 1.5, features exclusion of a portion of types - who are particularly risk averse and have very high costs of effort - with strictly positive measure whenever  $\omega > 0$ .

---

<sup>8</sup>With a slight abuse of notation,  $e(\mathbf{t}) = e(\alpha(\mathbf{t}), \mathbf{t})$  and  $\sigma(\mathbf{t}) = \sigma(\alpha(\mathbf{t}), \mathbf{t})$  correspond to the "recommended" choices introduced in the definition of contracts above.

Before we dive into the analysis of the incentive problems, we shall briefly discuss the case where the principal can observe the agent's type and his choices. If the principal is perfectly informed about the agent's preferences and choices, then there is no need to use the share  $\alpha$  to control incentives. Hence,  $\alpha$  is set so as to induce an optimal allocation of risk between agent and principal, so  $\alpha^*(\mathbf{t}) = 0$  for all  $\mathbf{t}$ . Since the principal is indifferent towards risk, and the efficient frontier is increasing in  $\sigma$ , he will prefer for any given effort the maximum volatility, so  $\sigma^*(\mathbf{t}) = \bar{\sigma}$  for all  $\mathbf{t}$ . Finally, the optimal level of effort satisfies the first-order condition  $\mu_e(e(\alpha^*(\mathbf{t}), \mathbf{t}), \bar{\sigma}) = c$ . Notice that both  $\alpha^*$  and  $\sigma^*$  are independent of the agent's preference parameters. The optimal level of the mean is decreasing in  $c$ , as effort is decreasing in  $c$  under complete information.

## 1.4 The Problem of Pure Moral Hazard

Even though we can characterize the solution to our model - in the case of commonly known characteristics - for general functions  $\mu(e, \sigma)$ , clear-cut comparative statics predictions require quite a lot more structure. Thus, to make progress we assume from now on that

$$\mu(e, \sigma) = e^\lambda \sigma^\delta. \tag{1.9}$$

It is easy to verify that the agent's problem of choosing  $e$  and  $\sigma$  for given contract is jointly concave in the choice variables if  $0 \leq \lambda, \delta \leq 1$  and  $\lambda + \frac{\delta}{2} \leq 1$ ,<sup>9</sup> so we impose these restrictions to make the agent's problem well behaved. As we discuss shortly, to make the principal's problem well behaved (that is concave in  $\alpha$ ), we assume on top of this that  $\lambda \leq .5$  and  $\delta \leq 2\lambda$ .

To solve our problem in the most reader friendly way, we proceed as follows. We demonstrate the important features of the solution for interior volatility choices in the main text. We provide the details of the solution in the appendix A.1 along side with a discussion for which parameter values the solution is indeed interior.

It is useful to ease notation defining some statistics of the model parameters. The details are not interesting in any way but are provided for completeness in Definition A.1 in the appendix A.1. Let  $\eta \equiv \eta(\delta, \lambda)$ ,  $\Delta \equiv \Delta(\delta, \lambda)$ , and  $\Gamma \equiv \Gamma(\delta, \lambda)$  denote functions of the parameters only, and let  $\theta(\mathbf{t}) \equiv a^{\frac{-\delta}{2(1-\lambda)-\delta}} \cdot c^{\frac{-2\lambda}{2(1-\lambda)-\delta}}$ . Building on this notation, we can write

---

<sup>9</sup>The factor .5 stems from the fact that the cost of risk bearing is quadratic in  $\sigma$ . Switching variables from standard deviation to variance in (1.9) gives rise to the standard restriction that the sum of exponents be smaller than unity.

the mean of the return distribution induced by an agent with characteristics  $\mathbf{t}$  as

$$\mu(e(\alpha, \mathbf{t}), \sigma(\alpha, \mathbf{t})) = \Delta\theta(\mathbf{t}) \alpha^{\eta-1}; \quad (1.10)$$

the agent's cost of risk bearing and effort as

$$a \frac{\alpha^2}{2} \sigma(\alpha, \mathbf{t})^2 + ce(\alpha, \mathbf{t}) = \Gamma\theta(\mathbf{t}) \alpha^\eta;$$

and the agent's indirect certainty equivalent level of wealth as

$$w_A(\beta, \alpha, e(\alpha, \mathbf{t}), \sigma(\alpha, \mathbf{t}), a) - ce(\alpha, \mathbf{t}) = \beta + \alpha\Delta\theta(\mathbf{t}) \alpha^{\eta-1} - \Gamma\theta(\mathbf{t}) \alpha^\eta. \quad (1.11)$$

### 1.4.1 Optimal Contracts

It is now straightforward to solve the principal's problem. Clearly, when the agent's characteristics are known, the participation constraint has to be binding for each  $\mathbf{t}$ . Notice, that the indirect certainty equivalent level of wealth - when the agent chooses effort and volatility optimally from his perspective - depends only on  $\theta = \theta(\mathbf{t})$ , a unidimensional statistic of  $\mathbf{t}$ . This simplifies the model dramatically. Imposing (1.7) for each  $\theta$  and substituting into the principal's objective, we obtain the following unconstrained problem

$$\max_{\alpha} \theta(\mathbf{t}) (\Delta\alpha^{\eta-1} - \Gamma\alpha^\eta). \quad (1.12)$$

It is easy to verify that problem (1.12) is increasing in  $\alpha$  for  $\alpha = 0$  and concave in  $\alpha$  for  $\eta \in (1, 2)$ , or equivalently for  $\lambda \leq .5$  and  $\delta \leq 2\lambda$ , which is precisely the reason we impose this restriction. In this case, we can characterize the solution by the pointwise first-order conditions

$$\alpha^* = \underline{\alpha} \equiv \frac{\eta - 1}{\eta} \frac{\Delta}{\Gamma}. \quad (1.13)$$

The optimal share  $\alpha^*$  is independent of the agent's preference parameters as long as the implied volatility choice is interior. This is due to the Cobb-Douglas technology. Since the agent's volatility choice is the higher the less risk averse the agent is and the smaller his marginal cost of effort is, the volatility choice is indeed interior for relatively high values of the agent's preference parameters. Let  $\mathbf{T}_I$  denote the (closed) set of parameters giving rise to an interior solution. Depending on the support of the agent's preference parameters, there is

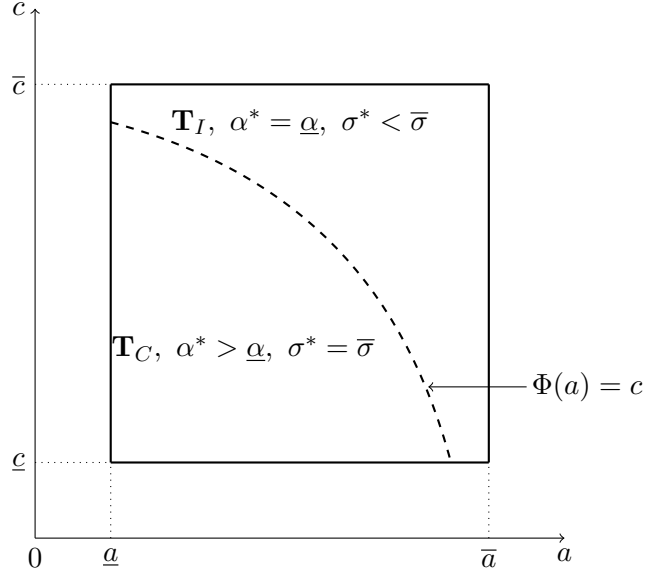


Figure 1.2: Interior versus corner solutions

necessarily a set of types for which the upper bound on volatility is a binding constraint. Let  $\mathbf{T}_C = \mathbf{T} \setminus \mathbf{T}_I$  denote this (open) set.  $\mathbf{T}_C$  is nonempty if some agents are close to risk neutral and/or have very low cost of effort. To capture this case, we assume that the lower bounds  $\underline{a}$  and  $\underline{c}$  are sufficiently low. In this case, there exists a strictly decreasing function  $\Phi(a)$  that separates the sets  $\mathbf{T}_C$  and  $\mathbf{T}_I$ , depicted in figure 1.2 below.

We can now characterize the overall solution to the contracting problem.

**Proposition 1.1** *i) There exists a strictly decreasing function  $\Phi(a)$ ,*

*such that  $\mathbf{T}_C \equiv \{\mathbf{t} : c < \Phi(a)\}$ . For  $\underline{a}$  and  $\underline{c}$  sufficiently small,  $\mathbf{T}_C$  is nonempty.*

*ii) For relatively risk averse agents with a relatively high cost of effort (formally, for  $\mathbf{t} \in \mathbf{T}_I$ ), the optimal share  $\alpha^*(\mathbf{t})$  is independent of  $\mathbf{t}$  and given by (1.13). The optimal choice of volatility and the expected level of profits are both decreasing in  $\mathbf{t}$ .*

*iii) For relatively risk tolerant agents with a relatively low cost of effort (formally, for  $\mathbf{t} \in \mathbf{T}_C$ ), the optimal share  $\alpha^*(\mathbf{t})$  is decreasing in  $\mathbf{t}$ . The optimal choice of volatility is  $\sigma^*(\mathbf{t}) = \bar{\sigma}$ . The expectation of profits is decreasing in  $\mathbf{t}$ .*

*iv) For any  $\mathbf{t}, \mathbf{t}'$  such  $\mathbf{t} \in \mathbf{T}_C$  and  $\mathbf{t}' \in \mathbf{T}_I$ , we have  $\alpha^*(\mathbf{t}) > \alpha^*(\mathbf{t}')$ .*

**Proof:** See in appendix A.1. ■

The economics is straightforward. An efficient allocation of risks would require that  $\alpha$

be set equal to zero for all  $\mathbf{t}$ . While a riskless contract would induce the agent to choose the optimal volatility from the principal's perspective, that is  $\sigma = \bar{\sigma}$ , it would not give the agent any incentive to exert effort. Hence,  $\alpha$  is set too high relative to the first-best. As a result, there is a strictly positive cost of risk bearing which is increasing in  $a$ . Since the agent's participation constraint always holds as an equality, it is the principal who bears the cost of this inefficiency. The higher is  $a$ , the more costly it becomes to convince the agent to participate for a given share of profits  $\alpha$ . Hence, the principal weakly reduces  $\alpha$  as  $a$  is increased. The agent, on the other hand, can reduce the cost of risk bearing by changing the volatility of the project. Hence, the agent (weakly) reduces the volatility of the project as he becomes more risk averse.

Similarly, when  $c$  increases, any given level of effort becomes more costly to implement. Hence, the principal finds it optimal to reduce incentives for effort when  $c$  increases, so  $\alpha$  is reduced. Since volatility and effort are complements along the efficient frontier, the agent has less of an incentive to engage in risk taking. Hence, the optimal volatility is reduced as well.

#### 1.4.2 Covariance of Contracts and Moments of the Profit Distribution

Inspired by Holmström & Milgrom (1994), we build our comparative statics predictions on the concept of associated random variables. Recall from Esary, Proschan & Walkup (1967) that random variables  $\mathbf{t}$  are associated if

$$COV(x(\mathbf{t}), y(\mathbf{t})) = \mathbb{E}[x(\mathbf{t})y(\mathbf{t})] - \mathbb{E}[x(\mathbf{t})]\mathbb{E}[y(\mathbf{t})] \geq 0$$

for all non-decreasing functions  $x(\mathbf{t})$  and  $y(\mathbf{t})$  (that is, functions that are non-decreasing in each of the arguments) for which  $\mathbb{E}[x(\mathbf{t})y(\mathbf{t})]$ ,  $\mathbb{E}[x(\mathbf{t})]$ , and  $\mathbb{E}[y(\mathbf{t})]$  exist.<sup>10</sup> Notice that the functions  $\alpha^*(\mathbf{t})$ ,  $\mu^*(\mathbf{t})$ , and  $\sigma^*(\mathbf{t})$  described in proposition 1.1 are comonotone. Before  $\mathbf{t}$  is realized, the values these functions take are random. Let  $\tilde{\alpha}^*$ ,  $\tilde{\sigma}^*$ , and  $\tilde{\mu}^*$  directly denote these random variables.

**Proposition 1.2** *i) The covariance of  $\tilde{\alpha}^*$  and  $\tilde{\sigma}^*$  is strictly positive.*

*ii) If  $\mathbf{t}$  is associated, then the covariance of  $\tilde{\alpha}^*$  and  $\tilde{\mu}^*$  and the covariance of  $\tilde{\mu}^*$  and  $\tilde{\sigma}^*$  are nonnegative.*

*iii) If either managerial risk aversion or his/her cost of effort can be controlled for, then the*

---

<sup>10</sup>For the relationship between association and other concepts of dependence (see Esary & Proschan 1972).

covariance of  $\tilde{\alpha}^*$  and  $\tilde{\mu}^*$  and the covariance of  $\tilde{\mu}^*$  and  $\tilde{\sigma}^*$  are both strictly positive.

**Proof:** See in appendix A.1. ■

Part i) follows from the fact that the functions  $\alpha^*(\mathbf{t})$  and  $\sigma^*(\mathbf{t})$  are comonotone and moreover that one function is strictly decreasing exactly in the region where the other is constant. Calculating the covariance by separating these regions yields the strictly positive result. Part ii) follows directly from the association property, because  $\alpha^*(\mathbf{t})$ ,  $\mu^*(\mathbf{t})$ , and  $\sigma^*(\mathbf{t})$  are all monotonic in  $\mathbf{t}$ . Finally, part iii) follows from the association property, the fact that one random variable is always associated, and that one can rewrite  $\tilde{\alpha}^*$  and  $\tilde{\sigma}^*$  as increasing functions of  $\tilde{\mu}^*$ .

The predictions of the pure moral hazard model when all moments of the profit distribution are endogenous are remarkably unambiguous: provided that the parameters in the agent's payoff function are positively correlated in the sense of association, the solution to the incentive problem and the induced moments reflect this positive correlation. This is in remarkably stark contrast to the predictions of the exact same model when the variance is taken as exogenous.<sup>11</sup> Intuitively, effort and risk are complements in the agent's problem, so both choices tend to increase the stronger are the incentives the agent faces. On the other hand, the principal offers steeper incentives to agents that are easier to incentivize.

## 1.5 The Case of Combined Adverse Selection and Moral Hazard

We now analyze the full problem, where the agent has private information about his preference parameters. In this case we obtain a rich set of comparative statics predictions also for the case where the feasibility constraint on the volatility is never binding. For convenience, we focus on this case.

Building on the analysis of the pure moral hazard case, we know that the agent's certainty equivalent level of wealth depends on the underlying parameters only through the statistic  $\theta$ . Therefore, it is clear that there must necessarily be bunching of types  $\mathbf{t}$  with the same level of  $\theta$ . From, (1.11) the agent's certainty equivalent level of wealth for any given announced type

---

<sup>11</sup>The standard trade-off between risk and incentives arises in the parameter set that gives rise to a corner solution,  $\mathbf{T}_C$ , when the upper bound on volatility,  $\bar{\sigma}$ , increases.



$\hat{\theta}$  is

$$\beta(\hat{\theta}) + (\Delta - \Gamma) \theta \alpha(\hat{\theta})^\eta.$$

Note that the cross derivative of this expression with respect to  $\alpha$  and  $\theta$  is positive, so the single crossing condition holds.<sup>12</sup> Moreover, the agent's indirect utility is the higher the higher is  $\theta$ .

A crucial difference between the present problem and the pure moral hazard problem is that the level of the agent's outside option matters quite a bit; the solution for the case where the agent's outside option,  $\omega$ , satisfies  $\omega > 0$  is qualitatively different from the case where  $\omega = 0$ , because the principal finds it optimal to exclude some types. Since types with higher  $\theta$  derive higher utility from participating, the principal excludes types with a low level of  $\theta$ .

Building on these insights, we can write the principal's problem formally as follows:

$$\begin{aligned} & \max_{\alpha(\cdot), \beta(\cdot), \theta_m} \int_{\theta_m}^{\bar{\theta}} \left\{ -\beta(\theta) + (1 - \alpha(\theta)) \theta \Delta \alpha(\theta)^{\eta-1} \right\} dF(\theta) \\ & \quad \text{s.t.} \\ & \beta(\theta) + (\Delta - \Gamma) \theta \alpha(\theta)^\eta \geq \beta(\hat{\theta}) + (\Delta - \Gamma) \theta \alpha(\hat{\theta})^\eta \text{ for all } \theta, \hat{\theta} \geq \theta_m \\ & \beta(\theta) + (\Delta - \Gamma) \theta \alpha(\theta)^\eta \geq \omega \text{ for all } \theta \geq \theta_m \\ & \max_{\hat{\theta}} \beta(\hat{\theta}) + (\Delta - \Gamma) \theta \alpha(\hat{\theta})^\eta \leq \omega \text{ for all } \theta < \theta_m. \end{aligned}$$

where  $F(\theta)$  is the cdf of the distribution of  $\theta$  and  $\theta_m$  is the marginal type  $\theta$  that is included; all types  $\theta < \theta_m$  are excluded.

The first step to solve this problem is to bring the incentive and participation constraint into a more tractable form. We call a pair of schedules implementable if they satisfy these two conditions.

**Lemma 1.1** *A pair of schedules  $\alpha(\theta)$  and  $\beta(\theta)$  is implementable if and only if*

$$\beta(\theta) = \int_{\theta_m}^{\theta} (\Delta - \Gamma) \alpha(z)^\eta dz - (\Delta - \Gamma) \theta \alpha(\theta)^\eta \text{ for all } \theta \geq \theta_m \quad (1.14)$$

---

<sup>12</sup>See Araujo et al. (2007) for an analysis of the case where the single crossing condition fails to hold. See Biais, Martimort & Rochet (2000) for a multidimensional model allowing for a similar reduction of the dimension of the incentive problem.

and  $\alpha(\theta)$  is nondecreasing in  $\theta$ . Exclusion of types  $\theta < \theta_m$  is incentive compatible if  $\alpha(\theta) = \beta(\theta) = 0$  for  $\theta < \theta_m$ .

**Proof:** The proof of the lemma is standard and therefore only sketched in the appendix A.2. ■

Only monotonic schedules  $\alpha(\cdot)$  can satisfy the incentive constraint. As monotonic schedules are differentiable almost everywhere, the agent's indirect utility function is differentiable almost everywhere. By the envelope theorem, the agent's utility changes with his preference statistic  $\theta$  at rate  $(\Delta - \Gamma)\alpha(\theta)^\eta \geq 0$ . Imposing the participation constraint for the marginal type  $\theta_m$  and integrating the changes in utility, we get (1.14). Finally, one shows that when contracts satisfy monotonicity of the schedule  $\alpha(\theta)$  and (1.14), then there is no profitable deviation for the agent. In particular, this argument implies also that deviations for a type  $\theta < \theta_m$  to any type  $\hat{\theta} \geq \theta_m$  would yield a level of utility that is strictly smaller than what the agent can get elsewhere,  $\omega$ .

Recall that  $\theta = \theta(\mathbf{t}) = a^{\frac{-\delta}{2(1-\lambda)-\delta}} \cdot c^{\frac{-2\lambda}{2(1-\lambda)-\delta}}$  is a statistic of the underlying parameters. Since  $\theta(\mathbf{t})$  is a function of random variables, we need to derive its distribution from the underlying distributions. The following lemma gathers the important features. For convenience, define  $r \equiv a^{\frac{-\delta}{2(1-\lambda)-\delta}}$  and  $s \equiv c^{\frac{-2\lambda}{2(1-\lambda)-\delta}}$ , respectively.

**Lemma 1.2** *The distribution of  $\theta$  is supported on a set  $[\underline{\theta}, \bar{\theta}]$ , where  $\underline{\theta} \equiv \underline{r}\underline{s}$  and  $\bar{\theta} \equiv \bar{r}\bar{s}$ . Moreover, let  $F(\theta)$  denote the cdf of  $\theta$  and  $f(\theta)$  denote the pdf. The density satisfies  $f(\underline{\theta}) = 0$  and  $f(\theta) > 0$  for  $\theta > \underline{\theta}$ . Moreover, provided that  $g(s|r)$ , the conditional density of  $s$  given  $r$  satisfies  $\frac{\partial}{\partial s}(sg(s|r)) \geq 0$ , the distribution of  $\theta$  satisfies  $\frac{\partial}{\partial \theta} \frac{1-F(\theta)}{\theta f(\theta)} \leq 0$  for  $\theta > \underline{\theta}$ .*

**Proof:** See in appendix A.2. ■

Note that the density of  $\theta$  goes to zero as  $\theta$  approaches the lower end of the type support. This is a well known property of this sort of problem and the driving force behind the exclusion result that we establish below, replicating Armstrong's (1996) observation for multidimensional screening problems more generally. As we discuss below in greater detail, the extent of exclusion in this particular contexts is simply a question of the level of the agent's outside option.

### 1.5.1 Optimal Contracts

It proves convenient to solve the principal's problem in two steps. In the first step, we take the choice of  $\theta_m$  as given and solve for optimal contracts for given  $\theta_m$ . In the second step, we address the exclusion problem. Types that are induced to opt out are offered a contract  $\alpha(\theta) = \beta(\theta) = 0$ . Substituting for  $\beta(\theta)$  from (1.14) into the principal's objective function and integrating by parts, we have

$$V(\theta_m) = \max_{\alpha(\cdot)} \int_{\theta_m}^{\bar{\theta}} \left\{ \left( \theta \Delta \alpha(\theta)^{\eta-1} - \Gamma \theta \alpha(\theta)^\eta \right) f(\theta) - (\Delta - \Gamma) \alpha(\theta)^\eta (1 - F(\theta)) \right\} d\theta - \omega(1 - F(\theta_m))$$

*s.t.*  $\alpha(\theta)$  nondecreasing in  $\theta$ .

The single crossing condition ensures that the participation constraint only binds at the low end of the support, in particular at  $\theta_m$ . For a given choice of  $\theta_m$ , the principal faces the standard efficiency versus rent extraction trade-off. On the one hand, the principal wishes to raise  $\alpha$  for each type so as to improve upon incentives for effort. On the other hand, the higher is  $\alpha$ , the higher are the rents the principal needs to give up to agents with a relatively high value of  $\theta$ . The optimal schedule  $\alpha(\theta)$  strikes a balance between these two motives. Under an appropriate regularity condition, the solution can be found by point-wise maximization under the integral. We state these results in the following proposition:

**Proposition 1.3** *Suppose that, for  $\theta > \underline{\theta}$ ,  $\frac{1-F(\theta)}{\theta f(\theta)}$  is non-increasing in  $\theta$ . Then, the optimal share schedule for  $\theta \geq \theta_m$  is given by*

$$\alpha(\theta) = \frac{\frac{\eta-1}{\eta} \Delta}{\Gamma + (\Delta - \Gamma) \frac{1-F(\theta)}{\theta f(\theta)}}.$$

*The optimal associated schedule  $\beta(\theta)$  is given by (1.14). In the limit as  $\theta_m \rightarrow \underline{\theta}$ , we have  $\lim_{\theta_m \rightarrow \underline{\theta}} \alpha(\theta_m) = 0$ .*

**Proof:** We omit a formal proof; the result follows straightforwardly from pointwise maximization. Moreover, it is easy to verify that the regularity condition implies that the solution is monotonic in  $\theta$ , so that we can indeed use pointwise maximization techniques. ■

The solution has the classical features. There is no distortion due to adverse selection for the agent with the highest parameter  $\theta$ ; that is,  $\alpha(\bar{\theta}) = \frac{\eta-1}{\eta} \frac{\Delta}{\Gamma}$ , corresponding exactly to the solution under pure moral hazard. For all  $\theta < \bar{\theta}$ , the share schedule is distorted downwards so as to extract rents from the agents with high parameters  $\theta$ . There is no rent at the bottom. Moreover, since the density of types  $\theta$  goes to zero at the low bound of the support, if such agents are offered a contract, then the shares they are offered become very small and go to zero as  $\theta_m \rightarrow \underline{\theta}$ . The reason is well understood from Armstrong (1996) and Rochet & Choné (1998). The density measures the weight given to the (constrained) efficiency motive in the principal's objective; on the other hand,  $1 - F(\theta)$  measures the weight given to the rent-extraction motive. Hence, at the low end of the support, the rent extraction motive becomes infinitely more important than the efficiency motive.

Consider now the optimal choice of types to include or to exclude, respectively. Using the first-order condition for the optimal  $\alpha(\theta)$ , the derivative of the principal's payoff with respect to  $\theta_m$  is

$$V'(\theta_m) = \left( \omega - \frac{\mu(\theta_m)}{\eta} \right) f(\theta_m),$$

where, with a slight abuse of notation,  $\mu(\theta)$  is short for the induced mean according to (1.10). The following results are now obvious:

**Proposition 1.4** *It is optimal to exclude a set of types with positive measure if and only if  $\omega > 0$ . The marginal type  $\theta_m^*$  is uniquely defined by the condition*

$$\omega\eta = \mu(\theta_m^*),$$

where  $\theta_m^*$  is the higher the higher is  $\omega$ . Moreover, the higher is  $\omega$ , the higher is the lowest incentive share that is offered,  $\alpha^*(\theta_m^*)$ , and the higher is  $\mathbb{E}[\alpha^*(\theta) | \theta \geq \theta_m^*]$ , the “average” observed incentive power of agents that are hired.

**Proof:** See in appendix A.2. ■

Since  $\alpha^*(\theta)$  goes to zero as  $\theta$  approaches the low end of the support, the expected profit generated by an agent of given type  $\theta$  goes to zero. Moreover, higher  $\theta$  types generate higher expected profits. Consequently, there is a uniquely defined marginal type  $\theta_m^*$  who generates exactly zero net surplus to the principal. Since only monotonic incentive schemes are incentive compatible, the positive implications of exclusion are as follows. The larger is

the set of excluded agents, that is the higher is  $\theta_m$ , the higher is the minimum level of  $\alpha$  that is observed in the cross section; in particular, the minimum exposure to risk is bounded away from zero so as to exclude some agents.<sup>13</sup>

## 1.5.2 Covariance of Contracts and Moments

We now turn to the comparative statics properties of the optimal contracting arrangement.

**Proposition 1.5** *i) With combined moral hazard and adverse selection, the covariance of  $\tilde{\alpha}^*$  and  $\tilde{\mu}^*$  is strictly positive whenever  $\alpha^*(\theta)$  is increasing in  $\theta$  on a set of positive measure.*

*ii) The covariance of  $\tilde{\alpha}^*$  and  $\tilde{\sigma}^*$  and of  $\tilde{\mu}^*$  and  $\tilde{\sigma}^*$ , respectively, is in general ambiguous. The covariance of  $\tilde{\alpha}^*$  and  $\tilde{\sigma}^*$  and of  $\tilde{\mu}^*$  and  $\tilde{\sigma}^*$ , respectively, is strictly positive if  $\mathbf{t}$  is associated and the distribution of  $\theta$  satisfies  $-\frac{\partial}{\partial \theta} \frac{1-F(\theta)}{f(\theta)} \leq \frac{\delta+2\lambda}{2(1-\lambda)-\delta}$  for  $\theta \geq \tilde{\theta}$ .*

**Proof:** See in appendix A.2. ■

Part i) of proposition 1.5 is due to the fact that  $\tilde{\alpha}^*$  and  $\tilde{\mu}^*$  are nondecreasing functions in  $\theta$ . Given  $\theta$  is unidimensional one can rewrite  $\tilde{\alpha}^*$  as a nondecreasing function of  $\tilde{\mu}^*$ . Since a scalar random variable is always associated, the result follows directly if the optimal  $\alpha$  is strictly monotonic on a set of positive measure. Part ii) states that, in general, the model loses its predictive power when it comes to the covariance of  $\tilde{\alpha}^*$  and  $\tilde{\sigma}^*$ . However, one can give simple sufficient conditions for a positive correlation between risk and incentives. The one given in proposition 1.5 ensures that the optimal profit share  $\alpha^*(\theta)$  does not change too fast as  $\theta$  changes, ensuring that the agent's optimal choice of  $\sigma$  becomes monotonic in the agent's underlying preference parameters. Since the defining property of associated random variables is precisely that the covariance of any monotonic functions of these random variables is positive, the conclusion follows immediately.

## 1.6 Attenuation

In our model, the moments of the profit distribution and the optimal contracts are endogenously determined as functions of the agent's preference parameters  $\mathbf{t} = (a, c)$ , that is, his degree of absolute risk aversion,  $a$ , and his marginal cost of effort,  $c$ . While we do not test our

---

<sup>13</sup>This should not be taken as a justification for high levels of manager compensation. The level is determined to a large extent by  $\omega$ , which is exogenous in the present model.

model directly, we now discuss how to bring it to the data and the consequences of neglecting the endogeneity of the variance.

If we are merely interested in contracts and risk, then a reduced form of our model is a system of equations for  $\alpha$  and  $\sigma$  (both normalized around their means) of the sort

$$\sigma = \gamma_1 a + \gamma_2 c + \varepsilon_1 \quad (1.15)$$

and

$$\alpha = \delta_1 a + \delta_2 c + \varepsilon_2, \quad (1.16)$$

where  $\varepsilon_1$  and  $\varepsilon_2$  are independent of each other and in particular independent of  $a$  and  $c$ .

What if the endogeneity of  $\sigma$  is neglected and instead  $\sigma$  is treated as a regressor for  $\alpha$ ? If there is some  $\beta_1$  such that  $\frac{\delta_1}{\gamma_1} = \frac{\delta_2}{\gamma_2} = \beta_1$ , then we can find another linear relation between  $\alpha$  and  $\sigma$  of the form

$$\alpha = \beta_1 \sigma + \varepsilon_3. \quad (1.17)$$

However, by implication of (1.15) and (1.16),  $\varepsilon_3$  is related to  $\varepsilon_1$  and  $\varepsilon_2$  according to

$$\varepsilon_3 = \varepsilon_2 - \beta_1 \varepsilon_1. \quad (1.18)$$

Using  $\mathbb{E}[\varepsilon_3] = 0$  and conditions (1.15) and (1.18), we find that

$$Cov(\sigma, \varepsilon_3) = \mathbb{E}[\varepsilon_3(\sigma - \mathbb{E}\sigma)] = \mathbb{E}[(\varepsilon_2 - \beta_1 \varepsilon_1)(\varepsilon_1)] = -\beta_1 Var(\varepsilon_1),$$

so that  $Cov(\sigma, \varepsilon_3) < (>) 0$  iff  $\beta_1 > (<) 0$ . Therefore, specification (1.17) fails to satisfy the assumptions of the linear regression model. As a consequence, the estimated value of  $\beta_1$  is biased towards zero, an effect that is known as attenuation (see e.g. Greene 1993); neglecting the endogeneity of  $\sigma$  biases the estimate of  $\beta_1$  towards zero. (The estimate is also inconsistent.)

The effects are slightly different but not more reassuring if  $\sigma$  is treated as a regressor alongside with controls  $a$  and  $c$ . Suppose we specify a linear model of the form

$$\alpha = \kappa_1 \sigma + \kappa_2 a + \kappa_3 c + \varepsilon_4 \quad (1.19)$$

in a situation where the true model is the system of equations (1.15) and (1.16). In fact, a form like (1.19) is obtained if we start from (1.16) and add  $\kappa_1$  times the difference between

the left and right side of (1.15). We obtain the following relation

$$\alpha = \kappa_1 \sigma + (\delta_1 - \kappa_1 \gamma_1) a + (\delta_2 - \kappa_2 \gamma_2) c + \varepsilon_2 - \kappa_1 \varepsilon_1,$$

which is indeed of the same form as (1.19) with  $\varepsilon_4 = \varepsilon_2 - \kappa_1 \varepsilon_1$ . Exactly the same attenuation problem as above arises if  $\delta_1 = \kappa_1 \gamma_1$  and  $\delta_2 = \kappa_2 \gamma_2$ . If  $\delta_1 \neq \kappa_1 \gamma_1$  or  $\delta_2 \neq \kappa_2 \gamma_2$ , then the direction of the bias is no longer clear, but the estimate of  $\kappa_1$  remains biased.

Summing up, neglecting the endogeneity of risk in simple regressions of contracts on measures of risk biases the estimates towards zero, irrespective of whether the estimated relation between the slope of incentive contracts and risk is positive or negative. In more sophisticated regressions that include risk as a regressor alongside with controls that effectively determine both the left- and the right-hand side of the regression equation, the direction of the bias is less clear; however, the estimation clearly remains biased also in these cases.

Whether, risk is exogenous or endogenous clearly depends on the context, so we cannot settle the question in a theoretical model. However, we simply point out, that attenuation makes it more likely to reject the hypothesis that incentives ( $\alpha$ ) depend positively (or negatively) on risk ( $\sigma$ ).

## 1.7 Conclusions

In this chapter, we analyze a model of managerial compensation with endogenous risk. Contracts serve a double purpose as providers of effort incentives and to guide the manager's project choices along an efficient frontier. The model offers a rich set of insights that have not been explored in such detail before. The resulting connection between risk and incentives depends on the underlying incentive problem. With pure moral hazard, a positive relation arises very naturally under general assumptions. With combined moral hazard and adverse selection, it is easy to find examples where the correlation between risk and incentives remains positive, but one can also construct cases where the covariance between risk and incentives is negative. However, we do not so much argue for a particular sign of this relation. The main point of the exercise is more that risk may be endogenous and to explore the implications of this variation. Empirically, endogeneity of risk gives rise to an attenuation problem resulting in estimates that are biased towards zero. We believe this may explain why a good part of the empirical studies on the subject produce relatively small (often statistically not significant)

relations between risk and incentives. We leave taking our model to the data directly to future work.

We have analyzed a particular incentive problem in this chapter where risk averse agents interact with risk neutral principals. As a consequence, our model cannot address excessive risk taking behavior. An incentive problem of this sort would arise, e.g., if both managers and principals are risk neutral and the manager gets some form of convex pay-scheme (e.g. through the use of options); similarly, excessive risk taking arises if managers and principals are risk neutral and firms suffer costs of financial distress - making the principal's payoff effectively concave in profits. Interestingly, even though the incentive problem differs substantially, the models may share the same comparative statics predictions that risk and incentives are positively related.







## Chapter 2

# Subgroup Deliberation and Voting

### 2.1 Introduction

Most committee decision making involves deliberation between heterogeneously informed individuals endowed with diverging preferences. Yet the interaction between the three aspects of information heterogeneity, preference heterogeneity and communication is non trivial. Heterogeneous information, in a common value setting, renders communication useful. Heterogeneity of preferences, on the other hand, makes communication difficult to achieve.

Committee communication, also called deliberation, always takes place according to some protocol which specifies a set of potential receivers and senders at every moment of time. Communication may be sequential or simultaneous. It may be entirely public, if messages are observed by everyone, or it may instead be semi-public, if communication is confined to Subgroups.

We examine two intuitive communication protocols in heterogeneous committees that vote under Unanimity: Plenary Deliberation and Subgroup Deliberation. Our aim is to rank these communication protocols w.r.t. simple Private Voting as well as among each other. We proceed in two main steps, by first isolating a set of equilibrium predictions for each protocol and then comparing these predictions as a means of comparing protocols.

The first step of our analysis is as follows. For each communication protocol as well as for Private voting, we restrict ourselves to a class of “simple” equilibria and call these respectively “Simple Subgroup Deliberation equilibria”, “Simple Plenary Deliberation equilibria” and “Simple No Deliberation Equilibria”. The restrictions on strategies embedded in the term “simple”

are mild in the case of Private Voting and in contrast significant in the case of Subgroup and Plenary Deliberation. Within the classes of equilibria considered, we furthermore only consider so called “reactive” equilibria, i.e. equilibria in which the same decision is not always made.

The second step of our analysis unfolds as follows. Having isolated a (non empty) set of equilibrium predictions for each of our protocols, we ask two specific questions. First, do there always exist reactive Simple Subgroup Deliberation and reactive Simple Plenary Deliberation equilibria that are Pareto improving w.r.t. any reactive Simple No Deliberation equilibrium? Secondly, does there always exist some reactive Simple Subgroup Deliberation equilibrium that is Pareto improving w.r.t any reactive Simple Plenary Deliberation equilibrium? Our answer to both questions is positive. The first result reveals that the two communication protocols dominate No Deliberation in a robust sense, given the mild restrictions imposed on strategies under Private Voting. Our second result shows that Subgroup Deliberation dominates Plenary Deliberation if one is willing to accept the significant restrictions that we impose on strategies under Plenary Deliberation. The latter form of dominance is thus admittedly significantly less general than the first form of dominance established. Modulo this important caveat, we thus obtain a complete ranking of the three voting mechanisms considered: Subgroup Deliberation dominates Plenary Deliberation which itself dominates Private Voting.

Among the plethora of potential communication protocols, we choose to focus on Plenary Deliberation and Subgroup Deliberation because we deem them intuitive and empirically relevant for the very reason that they are uncomplicated. The Plenary Deliberation protocol is equivalent to the common practice of straw votes: Each committee member simultaneously sends a public message chosen from a binary message space. Subgroup Deliberation restricts deliberation to homogeneous Subgroups. Examples of the latter protocol abound. In parliaments or parliamentary committees, party fellows often separately consult and reach a common stance before voting. Prior to faculty meetings, professors with related research agendas may meet separately. The key distinction between Plenary and Subgroup Deliberation resides in the a priori restriction that they place on information pooling. While Plenary Deliberation theoretically allows for a larger amount of information pooling than Subgroup Deliberation, our result is that Subgroup Deliberation however generates superior information sharing in equilibrium than Plenary Deliberation, when committees are heterogeneous.

In other words, our finding is that Subgroup Deliberation a posteriori generates more efficient information sharing than Plenary Deliberation for the very reason that it a priori restricts information sharing.

**Literature review** Early contributions in the literature on collective decision making and information aggregation focus on Private Voting and compare different voting rules. Seminal contributions such as Feddersen & Pesendorfer (1998), Gerardi (2000), and Duggan & Martinelli (2001) negatively single out Unanimity. Meirowitz (2002) adds a caveat to the above. The author examines a model featuring a continuum signal space as well as (at least nearly) perfectly informative signals and finds that full information equivalence obtains in the limit also for Unanimity.

Newer contributions add a stage of cheap talk communication prior to the vote. Gerardi & Yariv (2007) find that if one imposes no restriction on the communication protocol used, all non unanimous voting rules are equivalent in the sense that they induce the same set of equilibrium outcomes. Gerardi & Yariv (2007) contrasts with most of the remaining literature on cheap talk deliberation, which has instead examined specific protocols as well as simple equilibria. Most contributions have focused on the simultaneous Plenary Deliberation protocol and the truthful deliberation/sincere voting equilibrium (TS equilibrium). Coughlan (2000) shows that if preferences are known and substantially heterogeneous, the TS equilibrium does not exist. Austen-Smith & Feddersen (2006) show, within a generalized version of the classical Condorcet jury model, that uncertainty about preferences can render the TS equilibrium compatible with substantial heterogeneity, provided that the voting rule is not Unanimity. Meirowitz (2007), Van Weelden (2008), and Le Quement (2012) add further caveats to the analysis of Austen-Smith & Feddersen (2006). Finally, Deimen, Ketelaar & Le Quement (2014) show that if one considers a richer information structure featuring conditionally correlated signals, the TS equilibrium is compatible with a positive probability of ex post disagreement.

The question of the welfare properties of different protocols and equilibria has by and large been eluded. Clearly, in a homogeneous committee, the TS equilibrium implements the welfare maximizing decision rule, but little is known beyond this insight. Doraszelski, Gerardi & Squintani (2003) study a two persons setting with heterogeneous players who communicate simultaneously before voting under Unanimity. In equilibrium, information transmission is noisy, but communication is advantageous. Hummel (2012) identifies conditions under which

Subgroup Deliberation ensures no errors in asymptotically large and homogeneous committees. Wolinsky (2002) analyzes an expert game and shows that a Principal can sometimes gain by strategically grouping experts into optimally sized Subgroups that pool information before reporting to him.

This chapter complements existing literature on four aspects. First, it examines a little studied communication protocol, Subgroup Deliberation, that constitutes an alternative to Plenary Deliberation in heterogeneous committees in which types are publicly known. Second, it proposes a simple equilibrium scenario under Plenary Deliberation, for heterogeneous committees in which the TS equilibrium does not exist (so called minimally diverse committees (see Coughlan 2000)). Third, it provides a first attempt at a general clarification of the relative (Pareto) welfare properties of Private Voting, Subgroup and Plenary Deliberation. Finally, from a technical perspective, it introduces a simple method for the Pareto comparison of equilibria arising under different protocols in heterogeneous committees, which simply invokes a hypothetical sequence of best responses by different juror types.

The chapter is organized as follows. Section 2.2 introduces the basic jury model as well as the different communication protocols and equilibria that we consider. Section 2.3 provides a positive analysis of the equilibrium sets corresponding to the respective protocols under the imposed restrictions on strategy profiles. Section 2.4 compares the identified equilibria in terms of their Pareto welfare properties and thereby provides a tentative ranking of protocols. Section 2.5 concludes. Proofs are mostly relegated to appendixes B.1, B.2, and B.3.

## 2.2 The Model

### 2.2.1 Setup

Suppose a jury composed of  $n$  members. A defendant is being judged and is either guilty ( $G$ ) or innocent ( $I$ ) with equal prior probability. The jury must decide whether to convict ( $C$ ) or acquit ( $A$ ) him. Each juror casts a vote in favour of either conviction or acquittal. The voting rule is Unanimity: The defendant is convicted if and only if all jurors vote for conviction.

Each juror receives a single private signal prior to the vote. A signal  $s \in \{i, g\}$  indicates either guilt or innocence. A signal is “correct” with probability  $p \in (\frac{1}{2}, 1)$ , i.e.  $P(s = g |$

$G) = P(s = i | I) = p$ , while  $P(s = i | G) = P(s = g | I) = 1 - p$ . Juror signals are i.i.d. Let  $|g|$  denote the total number of  $g$ -signals received by the jury. The conditional probability  $P(G | |g| = k)$  that the defendant is guilty given  $|g| = k$  in an  $n$  persons jury is given as follows:

$$\beta(p, k, n) \equiv \frac{B(p, k, n)}{B(p, k, n) + B(1 - p, k, n)}, \text{ where } B(p, k, n) \equiv \binom{n}{k} p^k (1 - p)^{n-k}. \quad (2.1)$$

For  $j \in \{1, \dots, n\}$ , each jury member  $j$ 's preferences, are determined by a commonly known parameter  $q^j \in (0, 1)$ . A juror's payoff function is given as follows: Define  $U_j(C | I) = -q^j$  as the utility obtained by juror  $j$  when the defendant is convicted despite being innocent, and  $U_j(A | G) = -(1 - q^j)$  as the utility obtained when the defendant is acquitted but guilty. The utility related to remaining combinations of state and action (acquittal of an innocent or conviction of a guilty) is normalized to 0. Suppose a mechanism  $M$  yielding a probability  $P(C | I)$  of convicting an innocent defendant and a probability  $P(A | G)$  of acquitting a guilty defendant. The expected utility of juror  $j$  under mechanism  $M$  is given as follows:

$$U_j(M) \equiv -q^j P(C | I) P(I) - (1 - q^j) P(A | G) P(G). \quad (2.2)$$

Given this utility function, a juror  $j$  prefers conviction to acquittal whenever his posterior probability that the defendant is guilty exceeds  $q^j$ . The parameter  $q^j$  thus measures the juror's degree of aversion to wrongful conviction. The higher  $q^j$ , the more evidence of guilt is required for juror  $j$  to prefer conviction.

Juror preferences are heterogeneous and fall into two homogeneous categories. The jury contains  $n_D$  doves ( $D$ ) with preferences  $q_D$  and  $n_H$  hawks ( $H$ ) with preferences  $q_H$ , where  $q_H < q_D$  and  $n_D + n_H = n$ . We assume that at least one of the two preference types is present at least twice in the committee. We refer to the allocation of committee seats among preference types as the jury composition. For each  $j \in \{H, D\}$ , we use the notation  $-j = \{H, D\} \setminus j$ . For a given type  $j \in \{H, D\}$  and total number of signals  $\tilde{n}$ , the conviction threshold  $T_j^{\tilde{n}}$  is an integer number that satisfies the following:

$$\beta(p, T_j^{\tilde{n}} - 1, \tilde{n}) < q_j \leq \beta(p, T_j^{\tilde{n}}, \tilde{n}). \quad (2.3)$$

We make the following assumptions about preferences. First,

$$\mathbf{A.1:} \quad T_D^n - T_H^n \equiv m \geq 2.$$

In other words, in a putative equilibrium in which all  $n$  signals would be publicly revealed before the vote, at least two signal profiles would cause disagreement between the different juror types. The restriction is mild. Assuming  $m = 1$  typically imposes closely aligned preferences within the context of reasonably large committees in which many private signals are available. Second,

$$\mathbf{A.2:} \quad T_j^{n_j} \in \{1, \dots, n_j\}, \forall j \in \{H, D\}.$$

This means that if jurors of a given preference type  $j$  were to decide optimally on the basis of their  $n_j$  signals, they would sometimes acquit and sometimes convict. Finally,

$$\mathbf{A.3:} \quad q_D > \frac{1}{2}$$

This implies that a dove favours conviction only if the probability that the defendant is guilty exceeds  $\frac{1}{2}$ . This requirement matches the jury setting, where the "voir dire" selection process eliminates jurors that are excessively prone to convict. The assumption is used in proving our welfare results and we do not claim that it is necessary.

Throughout this chapter, we examine games exhibiting the following timing. In stage 0, jurors receive private signals. In stage 1, jurors communicate according to an exogenously fixed communication protocol. In stage 2, jurors simultaneously cast a vote. In stage 3, the defendant is convicted if and only if  $n$  conviction votes were cast.

## 2.2.2 Communication Protocols and Equilibria

We now introduce the three communication protocols that are the object of our analysis. No Deliberation (ND) simply specifies that no message is sent. Plenary Deliberation (PD) specifies that each juror simultaneously sends a message  $m \in \{i, g\}$  that is observed by all jurors. Subgroup Deliberation (SD) specifies that each juror simultaneously sends a message  $m \in \{i, g\}$  that is observed only by jurors of his preference type.

Protocols are orderable according to the physical restraints that they impose on communication. The first, No Deliberation, fully prohibits information sharing among jurors. The



second, Plenary Deliberation, potentially allows for full pooling of information among all jurors. The third, Subgroup Deliberation, prohibits communication between jurors of different preference types and only allows information pooling to take place within Subgroups of homogeneous jurors. Note that under Plenary as well as Subgroup Deliberation, we assume that communication is simultaneous, i.e. can be interpreted as simple straw votes preceding the actual vote. This is restrictive and must be distinguished from the free form communication considered in Gerardi & Yariv (2007).

We introduce a set of general definitions and restrictions on strategy profiles. A symmetric strategy profile specifies that jurors of the same preference type follow the same strategy. Monotonous strategies are s.t. information sets providing higher evidence of guilt are associated with a higher probability of voting for conviction. Throughout the analysis, we restrict ourselves to symmetric and monotonous strategies, in line with previous work on information aggregation and voting. We furthermore apply the following heuristic principle. For a given protocol, we ignore the possibility of mixing (in communication as well as in voting) as long as such a restriction does not leave us only with trivial equilibria in which the same decision (either  $C$  or  $A$ ) is always made. This is true of the PD and the SD cases. It is in contrast not true under ND and we thus consider the possibility of mixed voting under the latter protocol. We now present in detail the strategy profiles and equilibria that our analysis focuses on. Our focus is on perfect bayesian equilibria, which we simply call equilibria in what follows.

**No Deliberation** Under ND, jurors condition their votes exclusively on their own signal. We use the term “no deliberation strategy” instead of the standard term “private voting strategy” to describe the voting behavior of jurors under this protocol. A symmetric no deliberation strategy profile is characterized by a vector of mixing probabilities  $(\sigma_i^H, \sigma_g^H, \sigma_i^D, \sigma_g^D)$ , where  $\sigma_s^j$  denotes the probability that a single juror of type  $j$  votes for conviction given a signal  $s \in \{i, g\}$ . Let  $piv_j$  denote the event in which a given juror of preference type  $j$  is pivotal in the sense that the final decision changes with the juror’s vote. Let  $\gamma_G^j$  and  $\gamma_I^j$  denote the likelihood that a juror of preference type  $j$  votes for conviction given respectively state  $G$  or  $I$ . We have

$$\begin{aligned}\gamma_G^j &= p\sigma_g^j + (1-p)\sigma_i^j, \\ \gamma_I^j &= (1-p)\sigma_g^j + p\sigma_i^j.\end{aligned}$$

Define furthermore the indicator function  $Y(j, k)$  as follows. For  $j, k \in \{H, D\}$ ,  $Y(j, k) = 1$  if  $j = k$  while  $Y(j, k) = 0$  otherwise. Clearly, given the Unanimity rule,

$$P(G \mid s, piv_j) =$$

$$\frac{P(s \mid G) [\gamma_G^D]^{n_D - Y(j, D)} [\gamma_G^H]^{n_H - Y(j, H)}}{P(s \mid G) [\gamma_G^D]^{n_D - Y(j, D)} [\gamma_G^H]^{n_H - Y(j, H)} + P(s \mid I) [\gamma_I^D]^{n_D - Y(j, D)} [\gamma_I^H]^{n_H - Y(j, H)}}.$$

We call symmetric and monotonous no deliberation strategy profiles “simple ND profiles” (SND). If an SND profile is s.t. the defendant has a positive ex ante chance of both being acquitted or convicted, we call it a “reactive SND profile”. If an SND profile is s.t. the defendant is either always acquitted or always convicted, we call it a “non reactive SND profile”.

**Lemma 2.1** *Under the ND protocol, a reactive SND profile  $(\sigma_i^H, \sigma_g^H, \sigma_i^D, \sigma_g^D)$  constitutes an equilibrium iff,  $\forall j \in \{H, D\}, \forall s \in \{i, g\}$  :*

$$P(G \mid s, piv_j) = q_j, \text{ when } \sigma_s^j \in (0, 1), \quad (2.4)$$

$$P(G \mid s, piv_j) \leq q_j, \text{ when } \sigma_s^j = 0, \quad (2.5)$$

$$P(G \mid s, piv_j) \geq q_j, \text{ when } \sigma_s^j = 1. \quad (2.6)$$

**Proof:** The above conditions are standard (see e.g. Feddersen & Pesendorfer 1998) and their proof is therefore omitted. ■

Under the ND protocol, a reactive SND profile that constitutes an equilibrium is called a reactive SNDE.

**Plenary Deliberation** Under the PD protocol, consider first the strategy profile in which all jurors first truthfully reveal their signals while there is a threshold  $t \in \{1, \dots, n\}$  s.t. all jurors vote for conviction iff at least  $t$   $g$ -signals have been announced. We know from Coughlan (2000) that no such strategy profile constitutes an equilibrium of the game if  $m \geq 1$ . We instead examine a strategy profile that is given as follows. In Stage 1, jurors of type  $j$  truthfully reveal their signal while jurors of type  $-j$  simply always sends the message  $g$  and thus babble. In Stage 2, the voting decision of both juror types is conditioned on the number

of  $g$ -signals announced by type  $j$ . That is, there is a  $t_j \in \{0, 1, \dots, n_j, n_j + 1\}$  such that: 1) all jurors vote for conviction if at least  $t_j$   $g$ -signals have been announced by jurors of type  $j$  and 2) all jurors vote for acquittal otherwise. We call this strategy profile a “simple PD strategy profile” (SPD), thereby emphasizing the fact that one could envisage more complex strategy profiles under the PD protocol, for example involving noisy communication or mixed voting. We furthermore call an SPD profile a “reactive SPD profile” if  $t_j \in \{1, \dots, n_j\}$ , i.e. if jurors have a positive ex ante chance of unilaterally voting for both acquittal and conviction. If an SPD strategy profile is s.t. the defendant is either always acquitted or always convicted, we call it a “non reactive SPD strategy profile”.

Our restriction to pure strategies leaves us exclusively with equilibria in which doves truthtell while hawks babble. Truthtelling by doves appears natural given the allocation of power across types, which unambiguously favours doves. Given a profile of public information, if doves favour conviction, then hawks do so as well and will thus not veto such an outcome. If doves instead favour acquittal, they can furthermore always veto a conviction. In principle, doves can thus always get their way. The fact that hawks babble in the equilibria that we examine also appears quite natural in the light of this power allocation. As a matter of fact, we conjecture that there generally exists no symmetric and monotonic equilibrium in which an individual hawk is with positive probability pivotal at the communication stage. The argument behind this would be as follows. Given the preference misalignment assumed between doves and hawks ( $m > 1$ ), conditional on the event of being pivotal at the communication stage, a hawk favours conviction independently of his own signal. Consequently, if assumed to communicate informatively, a hawk will always favour announcing a  $g$ -signal.

**Lemma 2.2** *Under the PD protocol, a reactive SPD profile characterized by  $t_j \in \{1, \dots, n_j\}$  constitutes an equilibrium iff:*

$$\beta(p, t_j - 1, n_j) < q_j \leq \beta(p, t_j, n_j) \quad (2.7)$$

and

$$q_{-j} \leq \beta(p, t_{-j}, n_{-j} + 1). \quad (2.8)$$

**Proof:** The double inequality (2.7) is necessary and sufficient for a juror of type  $H$  not to have a strict incentive to deviate either at the communication or at the voting stage. The inequality (2.8) is necessary and sufficient to ensure that preference type  $-j$  is always willing

to vote for conviction whenever at least  $t_j$  guilty signals are announced by jurors of type  $j$ .

■

Under the PD protocol, a reactive SPD profile that constitutes an equilibrium is called a reactive SPDE. One may be uneasy with our ignoring the possibility of mixing at the voting stage. Our justification is purely practical: Including equilibria featuring mixed voting following truthtelling would be a daunting task for reasons that we explain in what follows. Recall that type  $j$  is the type that is truthtelling in the communication stage and consider an equilibrium featuring truthtelling followed by possibly mixed voting. Let  $(\theta_{-j}^i, \theta_{-j}^g)$  describe the (possibly mixed) voting strategy of type  $-j$ , where  $\theta_{-j}^s$  is the probability of voting  $C$  given signal  $s \in \{i, g\}$ . Symmetric mixed voting by jurors of type  $j$  requires indifference between decisions  $A$  and  $C$  at a given information set. This implies that given a voting strategy  $(\theta_{-j}^i, \theta_{-j}^g)$  of type  $-j$ , the mixed voting strategy of type  $j$  must be summarized by a vector  $(t_j, \theta_j)$  specifying the following voting behavior. When Subgroup  $j$  holds  $t_j$   $g$ -signals, each of its members votes  $C$  with probability  $\theta_j$ . When Subgroup  $j$  holds strictly more (less) than  $t_j$   $g$ -signals, all  $j$ -types convict (acquit). Furthermore, the conditional probability of guilt, conditional on  $t_j$   $g$ -signals in Subgroup  $j$  and on the assumption that all jurors of type  $-j$  convict, is equal to  $q_j$ . In order to characterize the set of equilibria featuring truthtelling followed by possibly mixed voting, one would thus have to identify an equilibrium vector given by  $(t_j, \theta_j, \theta_{-j}^i, \theta_{-j}^g)$ . This task is substantially more complicated than identifying a unique threshold  $t_j$  (equivalent to  $(t_j, 1, 1, 1)$ ) as we do. Furthermore, the increased complexity would carry over to the subsequent welfare exercise.

**Subgroup Deliberation** Under the SD protocol, we consider strategy profiles that are entirely characterized by a vector of thresholds  $t = (t_H, t_D)$ . In Stage 1, jurors simultaneously truthfully disclose their private signal to members of their Subgroup by sending a message identical to their signal. In Stage 2, all members of Subgroup  $j$  vote for conviction if the total number of guilty messages received among members of Subgroup  $j$  is weakly larger than  $t_j$ , and otherwise all vote for acquittal. We call this strategy profile a “simple SD profile” (SSD), thereby emphasizing the fact that one could construct more complex profiles under the SD protocol, for example involving noisy communication or mixing at the voting stage. We focus on SSD profiles that are such that the defendant has a positive ex ante chance of both being acquitted or convicted. We call such SSD profiles “reactive SSD profiles” and these come in

two subforms. A “type 2 reactive SSD profile” is a SSD profile in which  $t_j \in \{1, \dots, n_j\}$  for each  $j \in \{H, D\}$ . A “type 1 reactive SSD profile” is a reactive SSD profile in which one Subgroup  $j \in \{H, D\}$  adopts  $t_j = 0$ , while Subgroup  $-j$  adopts a threshold  $t_{-j} \in \{1, \dots, n_{-j}\}$ . If an SSD strategy profile is s.t. the defendant is either always acquitted or always convicted, we call it a “non reactive SSD strategy profile”.

We comment on key restrictions here. Given perfectly identical Subgroup preferences, focusing on outcomes featuring truthtelling appears natural. In contrast, one may be uneasy with our ignoring the possibility of mixing at the voting stage. Our justification is, as in the case of PD, purely practical: Including equilibria featuring mixed voting following truthtelling would be a daunting task. Symmetric mixed voting by jurors of type  $j$  requires indifference between decisions  $A$  and  $C$  at a given information set. This implies that given a strategy of type  $-j$  featuring truthtelling followed by (possibly mixed) voting, the mixed voting strategy of type  $j$  is summarized by a vector  $(t_j, \theta_j)$ , as in the case of mixed voting under PD described above. In order to characterize the set of equilibria featuring truthtelling followed by possibly mixed voting, one would thus have to identify an equilibrium vector given by  $(t_H, \theta_H, t_D, \theta_D)$ . This task is substantially more complicated than identifying a pair  $(t_H, t_D)$  (equivalent to  $(t_H, 1, t_D, 1)$ ) as we do. Furthermore, the increased complexity would carry over to the subsequent welfare exercise. More equilibria means more equilibria to compare, and mixed voting equilibria might not easily compare with each other or with pure voting equilibria. A final justification is the presumably limited impact of mixed voting on the set of implementable decision rules. When a Subgroup  $j$  is not excessively small, truthtelling in Subgroups implies a large array of revealed Subgroup signal profiles, out of which no more than one could induce randomized voting, as explained. When Subgroups are large, randomization in voting by a given preference type will thus only occur rarely in any given equilibrium and is thus arguably unlikely to heavily affect the type of implementable decision rules.

We now characterize conditions under which a given reactive SSD profile constitutes an equilibrium. Let  $|g|_j$  stand for the number of guilty signals held by Subgroup  $j$ . Let  $\left(|g|_j = t_j, |g|_{-j} \geq t_{-j}\right)$  denote the event in which Subgroup  $j$  holds exactly  $t_j$   $g$ -signals while Subgroup  $-j$  holds at least  $t_{-j}$   $g$ -signals.

**Lemma 2.3** *a) Under the SD protocol, a type 2 reactive SSD profile given by  $(t_H, t_D)$ , where*

$t_j \in \{1, \dots, n_j\} \forall j \in \{H, D\}$ , constitutes an equilibrium iff:

$$P\left(G \mid |g|_j = t_j - 1, |g|_{-j} \geq t_{-j}\right) < q_j \leq P\left(G \mid |g|_j = t_j, |g|_{-j} \geq t_{-j}\right). \quad (2.9)$$

b) Under the SD protocol, a type 1 reactive SSD profile given by  $(t_H, t_D)$ , where for some  $j \in \{H, D\}$ ,  $t_j \in \{1, \dots, n_j\}$  and  $t_{-j} = 0$ , constitutes an equilibrium iff (2.9) is true and

$$q_{-j} \leq P\left(G \mid |g|_{-j} = 0, |g|_j \geq t_j\right). \quad (2.10)$$

**Proof:** See in appendix B.1. ■

Under the SD protocol, a type 1 or type 2 reactive SSD profile that constitutes an equilibrium is called respectively a type 1 or type 2 reactive SSDE.

The idea behind reactive SSDEs is that each homogeneous Subgroup  $j$  votes as one person endowed with  $n_j$  signals. The SD protocol defines a sequential game in which individuals first communicate in Subgroups and then vote. We start with a discussion of Point a). The key insight is that condition (2.9) simultaneously ensures no strict deviation incentives both at the communication and at the voting stage. As to Point b), which characterizes type 1 reactive SSDEs, note that the behavior of Subgroup  $j$ , as specified in (2.9), is the same as if it were deciding alone and voting ex post optimally after fully pooling its information. Assuming that Subgroup  $-j$  convicts indeed provides no indication regarding the signal profile of the latter, as it always convicts. Subgroup  $-j$ , on the other hand, simply always convicts under the assumption that Subgroup  $j$  is convicting.

Our analysis unfolds in two steps. Section 2.3 provides a descriptive analysis of reactive SND, SPD and SSD equilibria. Section 2.4 analyzes the comparative welfare properties of reactive SSDEs, SPDEs and SNDEs.

## 2.3 Positive Analysis

**Lemma 2.4** *Under the ND protocol, a unique reactive SND profile constitutes an equilibrium. It is given by  $(\sigma_g^H = 1, \sigma_i^H = 1, \sigma_g^D = 1, \sigma_i^D = y)$ , where  $y \in (0, 1)$  if  $T_D^{n_D} < n_D$  and  $y = 0$  if  $T_D^{n_D} = n_D$ .*

**Proof:** See in appendix B.2. ■

The unique reactive SNDE, under our restrictions, is thus one in which hawks always convict, while doves vote as if they were an independent committee voting privately under Unanimity. The voting behavior of doves replicates the equilibrium characterized in Feddersen & Pesendorfer (1998). The key property of the unique reactive SNDE is that only the information of doves is aggregated, and typically imperfectly so, due to the fact that voting is private. As a final comment, note that our assumption that  $m > 1$  is key to eliminating a large amount of potential equilibrium scenarios under ND. When the doves are sufficiently biased towards acquittal (in relative terms), the assumption that all doves convict provides strong indication of guilt and unambiguously outweighs an individual hawk's information.

**Lemma 2.5** *Under the PD protocol, a unique reactive SPD profile constitutes an equilibrium. It is characterized by  $t_D = T_D^{n_D}$ .*

**Proof:** See in appendix B.2. ■

As already mentioned, it is intuitive that there exists an equilibrium in which doves publicly reveal their information, given that Unanimity voting effectively delegates decision power to them. This effective decision power of doves similarly explains why there is no reactive Simple Plenary Deliberation equilibrium in which hawks truthfully reveal their information. While the common feature of the unique reactive SNDE and SPDE is that hawks effectively delegate decision making to the doves, the difference between the two equilibria resides in the way doves aggregate their information. In the unique reactive SNDE, doves do not pool their information and thus always aggregate their information imperfectly if  $T_D^{n_D} < n_D$ . In the unique reactive SPDE, doves always fully pool their information, coordinate votes and aggregate their information optimally.

**Lemma 2.6** *Under the SD protocol:*

- a) *At least one reactive SSD profile constitutes an equilibrium.*
- b) *If there exist  $K > 1$  reactive SSDEs, then there exists a vector  $(t_H^1, t_D^1)$  s.t. the set of SSDEs is given by:*

$$(t_H^1, t_D^1), (t_H^1 - 1, t_D^1 + 1), \dots, (t_H^1 - K + 1, t_D^1 + K - 1). \quad (2.11)$$

**Proof:** See in appendix B.2. ■

Here again, there always exists an equilibrium satisfying our restrictions on strategies. In contrast to the sets of reactive SNDEs and reactive SPDEs, the set of reactive SSDEs may however contain more than one element. Point b) shows that if there exist several reactive SSDEs, these are orderable in terms of their degree of polarization. Among two reactive SSDEs, we say that the equilibrium with lower  $t_H$  and higher  $t_D$  is more “polarized”, because each of the Subgroups acts more in accordance with its own relative bias.

This concludes our descriptive equilibrium analysis, given our restrictions on strategy profiles. Having identified a set of equilibrium scenarios for each protocol, we may now proceed to a welfare comparison of the identified equilibria, aimed at producing a tentative ranking of the three considered protocols.

## 2.4 Normative Analysis

We say of an equilibrium that it is strongly Pareto dominant w.r.t. another equilibrium if both preference types obtain a strictly higher expected welfare in the first equilibrium. This subsection proceeds in three parts. First, Proposition 2.1 provides a Pareto welfare comparison of the unique reactive SPDE to the unique reactive SNDE. It establishes that the first equilibrium either strongly Pareto dominates the latter or is outcome equivalent to it. Second, Proposition 2.2 shows that when the set of reactive SSDEs is not a singleton, its elements are ordered in the strong Pareto sense. Third, Proposition 2.3 Pareto compares reactive SSDEs to the unique reactive SPDE. When the set of reactive SSDEs is not a singleton, the Pareto dominated equilibrium within this set either strongly Pareto dominates the unique reactive SPDE or is outcome equivalent to it. When the set of reactive SSDEs is a singleton, its unique element either strongly Pareto dominates the unique reactive SPDE or is outcome equivalent to it.

We add a comment on the interpretation of our theoretical exercise. Our reference to a jury setting may appear problematic because jury deliberations typically do not allow for Subgroup Deliberation. We see our analysis as a contribution to a normative debate aiming at potentially redesigning existing deliberation protocols in juries. In this perspective, considering new designs that are not in use seems legitimate. To the extent that one endorses our (admittedly restrictive) predictions for the different protocols, our welfare results would imply that members of a heterogeneous jury would unanimously agree to deliberate separately, if



given the choice between Plenary Deliberation and Subgroup Deliberation.

First, Jurors' ethnic or social background does appear to be a partial predictor of their preferences. Furthermore, the ethnic or social background of a person is at least imperfectly inferable from observable attributes (physical, verbal, psychological, etc).

**Proposition 2.1** *Reactive SPDE vs reactive SNDE.*

a) If  $T_D^{n_D} = n_D$ , the unique reactive SPDE is outcome equivalent to the unique reactive SNDE.

b) If  $T_D^{n_D} < n_D$ , the unique reactive SPDE is strongly Pareto dominant w.r.t the unique reactive NSDE.

**Proof:** See in appendix B.3. ■

As already mentioned, the unique reactive SNDE allows to optimally aggregate the information held by doves only if  $T_D^{n_D} = n_D$ , while the unique reactive SPDE always allows to achieve an optimal aggregation of the doves' information. This fact is reflected in the distinction between cases a) and b).

Our assumption that  $q_D > \frac{1}{2}$  is key to showing that the unique reactive SPDE strongly Pareto dominates the unique reactive SNDE if  $T_D^{n_D} < n_D$ . If  $q_D > \frac{1}{2}$ , a key aspect is that, maintaining the assumption of a unilateral conviction vote by hawks, transiting from private voting by doves (call this the "private" scenario) to an optimal aggregation of pooled signals by doves (call this the "pooled" scenario) leads to an increase in the ex ante probability of conviction and is thereby strictly beneficial to hawks. In the unique reactive SNDE, hawks indeed suffer from the doves' lack of willingness to convict. An adjustment in the doves' behavior that mitigates this reluctance without dramatically overshooting is thus naturally advantageous for hawks.

We now expand on the reason behind the fact that our condition requires a high enough  $q_D$ . As  $q_D$  increases, the probability of a unilateral conviction vote admittedly decreases under both scenarios ("private" and "pooled") considered above, but the key aspect is that this probability decreases faster under the first than under the second scenario. In the "private" scenario, a unilateral conviction vote by doves requires that every dove either receives a  $g$ -signal or, conditional on receiving an  $i$ -signal, votes for conviction, the latter event happening with probability  $y(p, q_D, n_D) \in (0, 1)$ . For very high values of  $q_D$ ,  $y(p, q_D, n_D)$  is however very

low and furthermore tends to 0 very fast as  $q_D$  tends to  $\beta(p, n_D - 1, n_D)$ . In contrast, as  $q_D$  increases and tends to  $\beta(p, n_D - 1, n_D)$ , the likelihood of a coordinated conviction vote by doves in the “pooling” scenario decreases slowly and without tending to 0. It is therefore quite intuitive that for  $q_D$  large enough, transiting from the “private” to the “pooling” scenario increases the likelihood of a unilateral conviction vote by doves.

Before going on to the final step of our normative analysis, which provides a comparison of reactive SSDEs to the unique reactive SPDE, we establish the preliminary result that the set of reactive SSDEs is fully orderable in the Pareto sense.

**Proposition 2.2 *Reactive SSDEs.***

*If  $(t_H, t_D), (t_H - 1, t_D + 1)$  are two reactive SSDEs, then  $(t_H, t_D)$  is strongly Pareto improving w.r.t.  $(t_H - 1, t_D + 1)$ .*

**Proof:** Consider two reactive SSDEs  $(t_H - 1, t_D + 1)$  and  $(t_H, t_D)$ . First, as proved in appendix B.3, transiting from  $(t_H - 1, t_D + 1)$  to  $(t_H - 1, t_D)$  is beneficial for the preference type  $H$  given our assumption that  $m > 1$ . Second, transiting from  $(t_H - 1, t_D)$  to  $(t_H, t_D)$  is also by definition beneficial to preference type  $H$ , given that  $t_H$  is type  $H$ 's best response to  $t_D$ . An equivalent argument shows that preference type  $D$  benefits from a transition from  $(t_H - 1, t_D + 1)$  to  $(t_H, t_D)$ . First, transiting from  $(t_H - 1, t_D + 1)$  to  $(t_H, t_D + 1)$  is beneficial for the preference type  $D$  given our assumption that  $m > 1$ . Second, going from  $(t_H, t_D + 1)$  to  $(t_H, t_D)$  is also by definition beneficial to preference type  $D$ , given that  $t_D$  is type  $D$ 's best response to  $t_H$ . ■

Proposition 2.2 shows that if there exist multiple reactive SSDEs, then the strongly Pareto dominant equilibrium within this set is easily described: it is that in which each preference type acts the least according to its own bias. In other words, it is the equilibrium in which the doves act harshest (have the lowest threshold  $t_D$ ) and the hawks act the most leniently (have the highest threshold  $t_H$ ). Reciprocally, the strongly Pareto dominated equilibrium within this set is the one in which preference types act the most in line with their relative bias. Summarizing, as one jumps from the one to the other adjacent equilibrium within the set of reactive SSDEs, the welfare of each type increases, the less that type acts in accordance with its relative bias.

We now finally compare reactive SSDEs with the unique reactive SPDE.

**Proposition 2.3** *Reactive SSDEs vs reactive SPDE.*

a) If  $q_H \leq P(G | |g|_H = 0, |g|_D \geq T_D^{n_D})$ , the type 1 reactive SSDE ( $t_H = 0, t_D = T_D^{n_D}$ ) exists and is outcome equivalent to the unique reactive SPDE. Any other reactive SSDE is strongly Pareto dominant w.r.t. the unique reactive SPDE.

b) If  $q_H > P(G | |g|_H = 0, |g|_D \geq T_D^{n_D})$ , any reactive SSDE is strongly Pareto dominant w.r.t. the unique reactive SPDE.

**Proof:** See in appendix B.3. ■

Proposition 2.3 builds on the following dynamic thought experiment: Start from the unique reactive SPDE, in which doves simply decide as if they were voting alone under Unanimity, fully pooling their information and optimally coordinating their votes according to the threshold  $T_D^{n_D}$ . Now, let hawks Subgroup Deliberate and optimally coordinate their votes under the assumption that doves convict, while doves continue to behave as in the unique reactive SPDE. There are now two possibilities, which are captured by respectively cases a) and b).

In case a), given that  $q_H \leq P(G | |g|_H = 0, |g|_D \geq T_D^{n_D})$ , hawks adopt a threshold  $t_H = 0$ . It follows that the type 1 reactive SSD profile ( $t_H = 0, t_D = T_D^{n_D}$ ) constitutes a reactive SSDE and is outcome equivalent to the unique reactive SPDE. In case b), given that  $q_H > P(G | |g|_H = 0, |g|_D \geq T_D^{n_D})$ , hawks instead adopt a threshold  $t_H > 0$ . This adjustment is by definition strictly improving for doves as well, as hawks become more lenient w.r.t. their previous voting behavior in the unique reactive SPDE.

We now expand on case b). The condition that  $q_H > P(G | |g|_H = 0, |g|_D \geq T_D^{n_D})$  means that the hawks' information is decision relevant in the sense that conditional on  $(|g|_H = 0, |g|_D \geq T_D^{n_D})$ , hawks favour an acquittal. Clearly, conditional on the information set  $(|g|_H = 0, |g|_D \geq T_D^{n_D})$ , the above condition implies that a dove would agree that an acquittal is optimal. Consequently, letting doves Subgroup Deliberate and coordinate votes according to  $T_D^{n_D}$ , both types gain if hawks now Subgroup Deliberate and coordinate votes according to some optimal threshold  $t_H > 0$  instead of always convicting. Now, let us consider a next round of adjustment: Let the doves optimally readjust their threshold in the light of the threshold  $t_H$  chosen by hawks in the previous round. It is clear that doves will choose  $t_D \leq T_D^{n_D}$ , so that this adjustment is at least weakly favourable to both preference types. This mutual adjustment process may be continued until a fixed point is reached. Such a fixed point exists

if there exists any reactive SSDE (and we know that there indeed exists one), and this fixed point corresponds to the most polarized reactive SSDE. Furthermore given that each step of the considered adjustment process is strongly Pareto improving, this reactive SSDE is strongly Pareto improving w.r.t. the unique reactive SPDE.

As a remark that applies to both cases a) and b) mentioned above, recall that if there exist several reactive SSDEs, we know from Proposition 2.2 that the most polarized reactive SSDE is strongly Pareto dominated by all remaining reactive SSDEs. It follows that if there are  $K > 1$  reactive SSDEs, then  $K - 1$  of these are a priori guaranteed to strongly Pareto dominate the unique reactive SPDE.

We now summarize our welfare comparison of the three protocols. Four cases can be distinguished.

The first and least interesting case corresponds to  $T_D^{n_D} = n_D$  and

$$q_H \leq P(G \mid |g|_H = 0, |g|_D \geq T_D^{n_D}). \quad (2.12)$$

Here, the unique reactive SPDE is outcome equivalent to the unique reactive SNDE and we furthermore cannot guarantee the existence of a reactive SSDE that strongly Pareto improves on the unique reactive SPDE. The only reactive SSDE that is guaranteed to exist is outcome equivalent to the unique reactive SNDE and SPDE.

The second case applies when  $T_D^{n_D} < n_D$  while (2.12) holds. Here, the unique reactive SPDE is strongly Pareto improving w.r.t. to the unique reactive SNDE and the only reactive SSDE of which we can guarantee the existence is outcome equivalent to the unique reactive SPDE.

The third case applies when  $T_D^{n_D} = n_D$  while (2.12) is reversed. Here, the unique reactive SPDE is outcome equivalent to the unique reactive SNDE and we know that there exists a reactive SSDE that strongly Pareto improves on the unique reactive SPDE.

The fourth and most interesting case applies when  $T_D^{n_D} < n_D$  while (2.12) is reversed. In this case, the unique reactive SPDE is strongly Pareto improving w.r.t. the unique reactive SNDE and we know that there exists a reactive SSDE that strongly Pareto improves on the unique reactive SPDE. We now summarize the intuition for this fourth case. One can think of the stepwise transition from ND to PD and then to SD in terms of two successive improvements. First, as compared to the unique reactive SNDE, the unique reactive SPDE

allows an improvement in the aggregation of the doves' information that is beneficial to both preference types. Secondly, as compared to the unique reactive SPDE, reactive SSDEs also allow to use the information held by the hawks, in a way that is advantageous to both preference types.

Given the above propositions, modulo our admittedly restrictive equilibrium selection under the PD and SD protocols, we have thus established a complete ranking of the three protocols considered: Subgroup Deliberation dominates Plenary Deliberation which itself dominates Private Voting. We wish to stress that the suboptimality of the ND protocol w.r.t. the remaining two protocols is a much more robust result than the dominance of SD over PD. Recall indeed that we impose very heavy restrictions on strategy profiles under PD and SD. Our ranking of SD and PD thus remains very tentative.

We close our analysis with two remarks on how our results potentially extend to more general settings. Our first remark concerns the condition  $q_D > \frac{1}{2}$  imposed throughout. As mentioned already, the condition is key to showing that the unique reactive SPDE strongly Pareto dominates the unique reactive SNDE if  $T_D^{n_D} < n_D$ . Now, assuming  $T_D^{n_D} < n_D$  and  $q_H > P(G \mid |g|_H = 0, |g|_D \geq T_D^{n_D})$ , we conjecture that one can construct examples in which  $q_D < \frac{1}{2}$  and the following holds true: The unique reactive SPDE is not Pareto improving w.r.t. the unique reactive SNDE, but some reactive SSDE however is. The rationale would be as follows: While the unique reactive SPDE is relatively unattractive in welfare terms, each step of the hypothetical adjustment process leading from the unique reactive SPDE to the most polarized reactive SSDE is Pareto improving and the set of reactive SSDEs is furthermore ordered in the Pareto sense.

## 2.5 Conclusion

We set out to compare three communication protocols characterized by different physical constraints on information pooling: PD, SD and ND. We identified simple conditions on juror preferences such that the following holds. First, the SD and PD protocols robustly dominate ND in the Pareto sense. The dominance of PD and SD w.r.t ND relies on the fact that the identified reactive SPDE and SSDE allow for a superior aggregation of the information held by doves, in a way that is also beneficial to hawks. Second, to the extent that one focuses on a restricted class of equilibria under PD, SD furthermore dominates PD. This second result

relies on the fact that the identified class of reactive SSDEs allows to also aggregate the information held by hawks.

Our analysis features a number of restrictions that future research should address. A truly robust comparison of PD and SD would need to characterize the whole set of reactive equilibria under each of the protocols, thus abandoning the restriction to monotonous, symmetric and pure strategies. It may be that PD and SD cannot be ranked in the Pareto sense. One also ought to consider other voting rules than Unanimity. In the case of SD and non unanimous voting rules, we conjecture that welfare dominant equilibria involve members of the same Subgroup voting asymmetrically. In such equilibria, the number of Subgroup members voting  $C$  would increase as a function of the number of  $g$ -signals held by the Subgroup. Another restriction of our analysis is the unrealistic assumption of only two preference types. Enlarging the set of preference types would however substantially complicate the analysis. One first direction to explore would be to assume that any juror's preference type is located within a neighbourhood of either of two reference values  $q_H$  or  $q_D$ . Finally, the binary information structure that we assume is restrictive. Our comparison of simple protocols ought to be repeated in a setting featuring continuous signals in order to evaluate whether our results still hold in such a more natural and versatile environment.







## Chapter 3

# Carbon Taxation under Asymmetric Information over Fossil-fuel Reserves

### 3.1 Introduction

Despite the increasing importance of CO<sub>2</sub> regulation, the question of how Pigovian taxation interacts with resource owners' incentives to report their reserves has not been analyzed yet. We study how a social planner being constrained to delayed regulation should define a carbon tax when the size of the fossil-fuel reserves is private information. We argue that the anticipation of CO<sub>2</sub> regulation may be a motive for over-reporting. Contrary to expectations, we demonstrate that also the regulator can profit from over-reporting, highlighting an “information curse” phenomenon.

The question of when the world will run out of fossil-fuels has received considerable attention in recent years. However, estimates of fossil-fuel reserves are subject to large errors, and contain different private information<sup>1</sup>. The main reasons for these errors are i) the lack of competition in fossil-fuel markets; ii) the spatial concentration of reserves; iii) the nature

---

<sup>1</sup>Oil owners in particular have private information on the size of their reserves. The International Energy Agency states: “Definitions of reserves and resources, and the methodologies for estimating them, vary considerably around the world, leading to confusion and inconsistencies. In addition, there is often a lack of transparency in the way reserves are reported: many national oil companies in both OPEC and non-OPEC countries do not use external auditors of reserves and do not publish detailed results.” (IEA 2010)

of the reserves assessment process that requires exploration permits, high technical skills, and dedicated capital; and iv) the absence of a worldwide control over resource assessments. The existing empirical studies (e.g. Bentley 2002, Laherrere 2013) have found that resource owners tend to over-report their reserves. This is typically explained by their incentives to discourage investments in R&D in backstop energy sources (Saure 2010), or to obtain larger production quotas in cartel organizations,<sup>2</sup> such as OPEC.<sup>3</sup>

According to the International Energy Agency (IEA 2011) CO<sub>2</sub> from energy production accounts for 65% of greenhouse gas emissions in 2009 and 81% of the world energy supply is based on fossil-fuels. The Intergovernmental Panel on Climate Change (IPCC 2013) states that CO<sub>2</sub> emissions are the largest cause of global warming. Plus, the IEA (IEA 2012) predicts that dangerous climate change cannot be prevented if fossil-fuel combustion won't be subject to restrictions. Even if the prevalent public opinion is that policy makers will not really act on these information, the chief economics commentator at the Financial Times, Martin Wolf,<sup>4</sup> advises fossil-fuel investors against totally dismissing the positive probability of a future legally binding obligation to prevent “dangerous anthropogenic (i.e., human-induced) climate change”.

This chapter is related to different strands of the literature. One strand is the literature on optimal pollution control of a stock of pollutants emerging from the use of an exhaustible resource where all information is publicly known. This has received considerable attention in a variety of settings based on the Hotelling model (Hotelling 1931): i) Some papers set a carbon ceiling on the CO<sub>2</sub> stock: e.g. Chakravorty, Magné & Moreaux (2006) consider this setting with one fossil-fuel and Chakravorty, Moreaux & Tidball (2008) deal with two fossil-fuels; ii) Others introduce a continuous increasing damage function: in this setup e.g. Ulph & Ulph (1994) and Tahvonen (1997) consider one fossil-fuel and Van der Ploeg & Withagen (2012) two fossil-fuels. The optimal carbon taxes in these papers are inverted U-shaped, due to natural dilution and the exhaustibility of the carbon-emitting resources. Ultimately, the fossil-fuel gets exhausted if natural dilution exists, the extraction costs are constant, and a relatively clean backstop is available.

---

<sup>2</sup>“Iraq, a founding member of OPEC, does not currently have a quota for crude production. Falah al-Amri, the head of the country’s State Oil Marketing Company, suggested that future quota calculations might have been a factor in the revision.” (<http://www.bloomberg.com/news/2010-10-04/iraq-lifts-oil-reserves-estimate-overtakes-iran-update1-.html>) seen 26. July 2014.

<sup>3</sup>Organization of the Petroleum Exporting Countries (OPEC).

<sup>4</sup>Martin Wolf 2014, “A climate fix would ruin investors”, Financial Times, 17. June, (<http://www.ft.com/cms/s/0/5a2356a4-f58e-11e3-afd3-00144feabdc0.html>) seen 25. June 2014.

In contrast to these studies, we consider an optimal tax in a world with private information on reserves. Also contrary to these papers, we use a discrete-time model with limited periods.

Strategic behaviors are excluded from the papers mentioned above. In contrast to these papers Wirl (1994) tackles a problem of an optimal pollution control with a strategic supplier of a polluting resource where all information is publicly known. Our model addresses strategic behaviors and shares some assumptions with Wirl's (1994) model.

Our work is also connected to the literature on asymmetric information in regulation problems. Jebjerg & Lando (1997) extend the Laffont-Tirole model of regulation under asymmetric information to pollution control, where the asymmetry of information is on pollution costs. Osmundsen (1998) examines the optimal taxation of the extraction of non-renewable natural resources under asymmetric information about the size of the reserves but does not introduce pollution problems. Both papers show that the optimal tax scheme is non-linear. The optimal tax in this chapter is non-linear, too.

A series of papers study the asymmetry of information over natural resources in delegation contracts. Gaudet, Lassere & Long (1995),<sup>5</sup>Osmundsen (1995) and Hung, Poudou & Thomas (2006) study the impact of information asymmetry of extraction costs on the extraction path and the tax revenues. Our work is the first one to address simultaneously: i) the optimal pollution control via a price instrument; and ii) private information on the reserves of polluting resources.

In our setup there is one cartel of resource owners<sup>6</sup> which sells its resources to consumers. The use of the exhaustible resource leads to emissions that accumulate in the atmosphere. The stock of pollutant creates environmental damage giving rise to a cumulative structure of the environmental damage. We assume in line with the literature that the environmental damage can be represented by an increasing and convex function of the carbon stock. A social planner whose objective is to maximize the expected weighted social welfare<sup>7</sup> taking

---

<sup>5</sup>See also (Ing 2012) that extends (Gaudet et al. 1995) to a framework where the capacity of the government to commit is limited. They do not consider the possibility of private information on the stock of resources.

<sup>6</sup>Oil markets have been modeled in various ways throughout the economic literature. The empirical literature has tried to determine the best market structure which explains the dynamics of the actual oil prices. Broadly speaking, there are two different approaches: one arguing that OPEC has the power to influence the market price (cartel, dominant firm model, target model) (Griffin 1985, Jones 1990, Dahl & Yucel 1991) and another one explaining prices either by property rights, scarcity or political events (Ezzati 1976, MacAvoy 1982, Verleger 1982, Loderer 1985). Our model is a cartel model, and hence belongs to the first class of models according to this classification.

<sup>7</sup>The weighted social welfare is a weighted combination of the consumers' utility and fossil fuel owner's surpluses

environmental damages into account implements a tax in a second period. Since no worldwide taxation of carbon emissions is implemented at the moment, we assume that the social planner cannot set a carbon tax in the first period. In other words, she is not able to tax the first period consumption, as she needs time to implement her optimal policy, and accelerating negotiations is too costly. This feature of the model is what we call “delayed regulation”. We assume that there is a consensus about the existence of a future carbon tax. This type of delayed regulation is plausible; it echoes modeling choices of the literature where the future carbon tax path is perfectly anticipated and the tax is globally increasing over a first segment of time (see Bauer, Hilaire & Bertram 2014). In our case, the first period tax is simply zero by assumption. The tax revenues are lump-sum redistributed to the consumers and hence enter the regulator’s utility function.

We consider a 3 period game.<sup>8</sup> The timing of the game is the following. In period 0, nature draws the size of the fossil-fuel reserves and the resource owner observes the size. In period 1, the fossil-fuel owner sets a price for fossil-fuels, the consumers demand a part of the reserves via the given demand function, and the social planner observes this amount. Given her observations she estimates the size of the fossil-fuel reserves, taking into account that the fossil-fuel owner tries to influence the endogenous carbon policy of period 2 in order to maximize his profits. Between period 1 and period 2 the regulator announces a tax which has to be paid by the consumers on their demand in period 2.<sup>9</sup> The tax is set to maximize the regulator’s expected utility of period 2, given her anticipated supply in period 2. Once her policy is announced, the social planner cannot change it. In period 2, the fossil-fuel owner again sets a price for fossil-fuels but this time taking the tax into account, the consumers demand a part of the remaining reserves given the demand function and pay the tax on their demand. In period 3, a clean backstop technology is available, in an infinite quantity at a very low cost, which eliminates the demand for fossil-fuels from that period on.

To understand the driving forces of the model we first consider the public information case in which the size of the fossil-fuel reserves is common knowledge. Due to the cumulative structure of the environmental damage, the social planner prefers a lesser supply in period 1 and a higher supply in period 2 compared to the fossil-fuel owner’s preferred supplies.

The social planner can use the tax as an aid to force the resource owner to supply the

---

<sup>8</sup>For a survey on dynamic games in the natural resources setup (see Long 2011).

<sup>9</sup>On the issue of permanent versus interim regulations (see Malik 1991).

amounts which are more in her interests, by punishing the fossil fuel owner with a high tax if he sells too much resources in period 1 and rewarding him with a low tax if the supply of period 1 is in her interests. Taking the social planner's strategy into account the fossil-fuel owner is going to trade off the effect of the social discount factor on his utility against the effect of the social planner's tax strategy on his utility when choosing his preferred supplies of period 1 and period 2.

The size of the reserves can be classified into two different groups by a threshold amount, i) if the reserves are higher than this threshold, the resources will not be exhausted in equilibrium. The corresponding equilibrium strategies are independent of the size of the reserves. The resource owner supplies unique optimal amounts of fossil-fuels in each of the two periods and at the same time, the regulator implements a unique optimal tax, ii) if, on the other hand, the reserves are lower than this threshold, the consumers' total demand exceeds the size of the reserves and thus the fossil-fuels will be exhausted in equilibrium. The associated equilibrium strategies are not independent of the size of the reserves. As in equilibrium the total supply is below the favored total demand, also the regulator prefers total exhaustion despite the environmental damages. Since the tax revenues enter the social planner's utility, it turns out that her best tax response is higher if the fossil-fuel owner has lower reserves left in period 2 than if he has higher reserves left in period 2. Rephrasing, she "discriminates" fossil-fuel owners of lower reserves against fossil-fuel owners of higher reserves, if they have supplied the same amount in period 1.

Let us now turn to the private information case. We assume that the fossil-fuel owner has either high reserves (high type) or low reserves (low type) and is privately informed about the size of the reserves. We find that the regulator's tax response if she does not know which type she is facing is a weighted average of the symmetric information tax responses of both types. Given this tax strategy and the fact that the regulator "discriminates" fossil-fuel owners of lower reserves against fossil-fuel owners of higher reserves, the high type does not have an incentive to mimic the low type. The low type, on the other hand, trades off his gain from facing a lower tax given that he mimics the high type in period 1 to his loss from deviating from his preferred supply of period 1. We identify conditions such that there always exists an equilibrium in which the high type supplies his symmetric information equilibrium quantities in each of the two periods. This equilibrium is a pooling equilibrium if the low type's reserves have the capacity to mimic the high type and it is a separating equilibrium if this is not

the case. The described equilibrium –pooling or separating– is ex ante preferred compared to any other existing equilibrium by both types of fossil-fuel owners. And it is the only equilibrium which survives the Cho-Kreps’ Intuitive Criterion. Moreover, if this equilibrium is a separating equilibrium we show that there does not exist any other equilibrium. However, if this equilibrium is a pooling equilibrium there might exist other equilibria. On top of this, we highlight that the social planner can ex ante prefer this equilibrium to the two symmetric information equilibria given each type. This, in a sense, means that the regulator can profit from asymmetric information compared to symmetric information.

The remainder of this chapter is organized as follows. Section 3.2 presents the model. Section 3.3 studies the case where the information of reserves is common knowledge. Section 3.4 presents the private information case. Section 3.5 does some welfare analysis. Last, Section 3.6 concludes. Proofs are mostly relegated to the appendices C.1, C.2, and C.3.

## 3.2 The Model

Our model shares some assumptions with Wirl’s (1994) model. In contrast to his work we analyze an asymmetric information game, and regulation –through a consumers’ quantity tax  $\tau$ – is only implemented in period 2.

**Basic properties of the model** Denote by  $Q \geq 0$  the size of the fossil-fuel reserves and by  $r > 0$  the social discount factor. Let the energy demand function be a time invariant linear function of the consumers’ fossil-fuel price  $p$ ,  $D(p) = 1 - p$ .

The fossil-fuel owner’s revenues per unit sold when the consumers face a quantity tax,  $\tau \in [0, 1]$ , is  $p - \tau$ . We assume that all fossil-fuel owner’s costs, like the extraction, exploration, and production costs, are zero. His profit,  $\pi$ , of selling  $D(p)$  amount of fossil-fuels at the price  $p$  and given a quantity tax,  $\tau$ , therefore writes

$$\pi(p, \tau) = (p - \tau)D(p) = (p - \tau)(1 - p) = (1 - D(p) - \tau) D(p) = \pi(D(p), \tau). \quad (3.1)$$

We consider a model in which fossil-fuels are sold in two periods and the social planner taxes the consumption of period 2 via a quantity tax,  $\tau$ . The amount sold in period 1 is denoted by  $Q_1$  and the amount sold in period 2 by  $Q_2$ .

Hence, the fossil-fuel owner's present value at period 1 is

$$\begin{aligned} U_F(Q_1, Q_2, \tau) &\equiv \pi(Q_1, 0) + \frac{1}{1+r} \pi(Q_2, \tau) \\ &= (1 - Q_1) Q_1 + \frac{1}{1+r} (1 - Q_2 - \tau) Q_2. \end{aligned} \quad (3.2)$$

And, the consumers' present value at period 1, assuming that their processing costs are zero, is

$$\begin{aligned} U_C(Q_1, Q_2, \tau) &\equiv \int_0^{Q_1} (1 - \tilde{q}) d\tilde{q} - (1 - Q_1) Q_1 + \frac{1}{1+r} \left( \int_0^{Q_2} (1 - \tilde{q}) d\tilde{q} - (1 - Q_2) Q_2 \right) \\ &= \frac{1}{2} Q_1^2 + \frac{1}{1+r} \frac{1}{2} Q_2^2. \end{aligned} \quad (3.3)$$

And finally, the present tax revenues at period 1 are

$$\Upsilon(Q_1, Q_2, \tau) \equiv \frac{1}{1+r} \tau Q_2. \quad (3.4)$$

We assume that the social planner's objective is to regulate carbon emissions in order to maximize a weighted social welfare taking external costs as environmental damage into account: i) The weighted social welfare is a weighted average of the consumers' and fossil-fuel owner's welfare. The tax revenues are lump-sum redistributed to the consumers and hence enter the regulator's utility. The weighted social welfare of and after period 3 is simply zero as we assume that a clean backstop freezes the demand for fossil-fuels from that period on. The weighted social welfare at period 1 can be stated as

$$W(Q_1, Q_2, \tau) \equiv \lambda U_F(Q_1, Q_2, \tau) + (1 - \lambda) [U_C(Q_1, Q_2, \tau) + \Upsilon(Q_1, Q_2, \tau)], \quad \lambda \in [0, \frac{1}{2}). \quad (3.5)$$

ii) The regulator's cost of carbon emissions, which reflects the environmental damage, is assumed to be a quadratic function,  $C(q) = \frac{d}{2} q^2$ , of the consumed amount of fossil-fuels,  $q \geq 0$ , where  $0 < d < \frac{2(1-\lambda)r}{1+r}$ . The damage coefficient,  $d$ , measures the impact of one unit burned fossil-fuels on the environment. Climate change is triggered by carbon exceeding a natural carbon stock that equals the initial carbon stock in the atmosphere. Hence, the regulator's dis-utility of the amount,  $Q_1$ , consumed in period 1 is  $-C(Q_1)$ . To the regulator's cost of emissions in following periods note that the environment is not only affected by the

carbon flow but in particular by the carbon stock. So, the dis-utility of carbon emissions has a negative effect on the social planner's utility in all future periods. Therefore, we assume that the social planner's dis-utility in and after period 2 is  $-C(Q_1 + Q_2)$ .<sup>10</sup>

The regulator's present utility at period 1 is hence

$$U_S(Q_1, Q_2, \tau) \equiv W(Q_1, Q_2, \tau) - C(Q_1) - \frac{1}{r}C(Q_1 + Q_2). \quad (3.6)$$

**Asymmetric information** We consider a signaling model with asymmetric information on the size of the fossil-fuel reserves. The reserves are assumed to be either low  $Q^L \in [0, \infty)$  or high  $Q^H \in (0, \infty)$ , with  $Q^L < Q^H$ . And they are distributed according to the commonly known probabilities  $p_r(Q^L) = \mu$ ,  $p_r(Q^H) = 1 - \mu$ , with  $0 < \mu < 1$ .

**The timing of the game**

- At  $t = 0$ , nature draws the fossil-fuel owner's type,  $Q$ , from the set  $\{Q^L, Q^H\}$ , with the above described probability distribution,  $p_r$ . The fossil-fuel owner observes his type.
- At  $t = 1$ , the fossil-fuel owner sets a price for fossil-fuels. Given this price, the consumers demand an amount  $0 \leq Q_1 \leq Q$  of the fossil-fuels, which satisfies their demand function.
- Between period 1 and period 2 after observing the amount,  $Q_1$ , the social planner forms a belief about the fossil-fuel owner's type  $\mu^*(Q^L | Q_1) = 1 - \mu^*(Q^H | Q_1)$ . Depending on  $Q_1$  and the belief,  $\mu^*$ , the regulator announces a tax,  $\tau(Q_1)$ , which has to be paid by the consumers<sup>11</sup> on their demand in period 2.
- At  $t = 2$ , the fossil-fuel owner sets a price for fossil-fuels,  $p$ , by taking the tax,  $\tau$ , into account. The consumers demand a part  $Q_2 \leq Q - Q_1$  of the remaining fossil-fuels given their demand function at the price  $p$ . The tax the consumers have to pay in period 2 is  $Q_2\tau$ . And the fossil-fuel owner's revenues per unit sold is the price shifted with respect to the tax level,  $p - \tau$ .
- At  $t = 3$ , a clean backstop technology is available in an infinite amount and at a very low cost, such that it destroys the demand for fossil-fuels in all following periods.

As the price of the fossil-fuels and the amount consumed have a one-to-one relationship via the demand function, for the remainder of this chapter, we view the fossil-fuel owner's strategy as deciding which amount of fossil-fuels to supply instead of setting a price for fossil-fuels.

<sup>10</sup>The dis-utility of environmental damage in all periods following period 2 is also  $-C(Q_1 + Q_2)$ , as we assume that in period 3 a clean backstop freezes the demand for fossil-fuels in all following periods.

<sup>11</sup>Instead claiming that the fossil-fuel owner pays the tax on his supply eventuates in the same results.



So, the regulator's objective between period 1 and period 2, of levying a quantity tax,  $\tau(Q_1)$ , is to maximize the weighted social surplus of period 2 minus the cost of environmental damage of and after period 2. The environmental damage is triggered by the supply,  $Q_1$ , of period 1 and the fossil-fuel owner's response  $Q_2(\tau)$  to the tax,  $\tau$ , in period 2. Hence, the regulator's problem between period 1 and period 2 is to maximize his utility of period 2,

$$U_S^{t=2}(Q_1, Q_2(\tau), \tau) \equiv \lambda[1 - Q_2(\tau)]Q_2(\tau) + \frac{1}{2}(1 - \lambda)Q_2(\tau)^2 + (1 - 2\lambda)Q_2(\tau)\tau - \frac{1+r}{r} \frac{d}{2}[Q_1 + Q_2(\tau)]^2, \quad (3.7)$$

with respect to  $\tau$ .

**The equilibrium conditions** A strategy profile,  $(Q_1^{L*}, Q_2^{L*}, Q_1^{H*}, Q_2^{H*}, \tau^{L*}, \tau^{H*})$ , and a belief system,  $\mu^*$ , constitute an equilibrium if and only if they satisfy:

The fossil-fuel owner's problem:

$$\{Q_1^{x*}, Q_2^{x*}\} = \operatorname{argmax}_{Q_1, Q_2} U_F(Q_1, Q_2, \tau^*(Q_1)) \quad (3.8)$$

s.t.

$$Q_1^{x*} + Q_2^{x*} \leq Q^x, \text{ for all } x \in \{L, H\}. \quad (3.9)$$

Additionally for all  $Q_1$  the tax  $\tau^*(Q_1)$  has to satisfy the regulator's problem:<sup>12</sup>

$$\tau^*(Q_1) = \operatorname{argmax}_{\tau} [\mu^*(Q^L | Q_1) U_S^{t=2}(Q_1, Q_2^{L*}(Q_1, \tau), \tau) + (1 - \mu^*(Q^L | Q_1)) U_S^{t=2}(Q_1, Q_2^{H*}(Q_1, \tau), \tau)], \quad (3.10)$$

where

$$Q_2^{x*}(Q_1, \tau) = \operatorname{argmax}_{Q_2} U_F^{t=2}(Q_2, \tau) \text{ s.t. } Q_2 \leq Q^x - Q_1, \text{ for all } x \in \{L, H\}, \quad (3.11)$$

with  $U_F^{t=2}(Q_2, \tau) \equiv (1 - Q_2 - \tau) Q_2$ .

---

<sup>12</sup>The equilibrium tax is  $\tau^{x*} = \tau^*(Q_1^{x*})$ .

### 3.3 The Case of Symmetric Information

Let the size of the fossil-fuel reserves be public information and either low,  $Q^L$ , or high,  $Q^H$ , with  $0 \leq Q^L < Q^H \leq Q_k^1 \equiv Q_k^1(d, r, \lambda)$  and  $x \in \{L, H\}$ .<sup>13</sup> We identify the Subgame Perfect Nash equilibrium strategies,  $Q_1^{x**}, Q_2^{x**}$  and  $\tau^{x**}$ , of this game. Under symmetric information, the general conditions given in (3.8), (3.9), (3.10), and (3.11) reduce to the following conditions. A strategy profile,  $(Q_1^{x**}, Q_2^{x**}, \tau^{x**})$ ,  $x \in \{L, H\}$ , constitutes an equilibrium if and only if it satisfies:

The fossil-fuel owner  $x$ 's problem:

$$\{Q_1^{x**}, Q_2^{x**}\} = \operatorname{argmax}_{Q_1, Q_2} U_F(Q_1, Q_2, \tau^{x**}(Q_1)) \quad (3.12)$$

s.t.

$$Q_1^{x**} + Q_2^{x**} \leq Q^x. \quad (3.13)$$

The regulator's problem: For any  $Q_1$  the tax  $\tau^{x**}(Q_1)$  has to satisfy

$$\tau^{x**}(Q_1) = \operatorname{argmax}_{\tau} U_S^{t=2}(Q_1, Q_2^{x**}(Q_1, \tau), \tau),^{14} \quad (3.14)$$

where

$$Q_2^{x**}(Q_1, \tau) = \operatorname{argmax}_{Q_2} U_F^{t=2}(Q_2, \tau) \quad \text{s.t.} \quad Q_2 \leq Q^x - Q_1, \quad x \in \{L, H\}. \quad (3.15)$$

Before we start to analyze the model let us first have a look at the essence and the technical properties of the threshold,  $Q_k^1$ , and its consequences for our analysis.

Note first, if sufficiently large fossil-fuel reserves exist the equilibrium problem can be solved without the quantity constrain (3.13). Indeed, solving the problem without the quantity constraint yields “optimal demands” of fossil-fuels in both periods,  $Q_1^{op}$  and  $Q_2^{op}$ , and an “optimal tax”,  $\tau^{op}$ . These strategies do not directly depend on the level of the reserves,  $Q^x$ ,  $x \in \{L, H\}$ , as the players' utilities do not directly depend on it. When the fossil-fuel reserves,  $Q^x$ , are sufficiently high the “optimal demands” of both periods can be supplied,

<sup>13</sup> $Q_k^1(d, r, \lambda)$  is a threshold such that if and only if the reserves,  $Q$ , satisfy,  $Q \geq Q_k^1$ , the equilibrium strategies are identical. The threshold is defined in more detail in the appendix C.1.

<sup>14</sup>The equilibrium tax is  $\tau^{x**} = \tau^{x**}(Q_1^{x**})$ .

$Q_1^{op} + Q_1^{op} \leq Q^x$ , and consequently the quantity constraint (3.13) is indeed redundant.

If the reserves instead satisfy,  $Q^x \leq Q_1^{op} + Q_1^{op}$ , the quantity constraint (3.13) is binding in equilibrium,  $Q_2^{x**} = Q^x - Q_1^{x**}$ , and the equilibrium strategies consequently depend on the level of the reserves.

The threshold,  $Q_k^1$ , is chosen in such a way that the described “optimal demands” can be just supplied. In other words, it is the sum of the “optimal demands”,  $Q_k^1 \equiv Q_1^{op} + Q_1^{op}$ . So, if the reserves,  $Q^x$ , are higher than this threshold,  $Q^x \geq Q_k^1$ , the equilibrium strategies are identical and equal the “optimal strategies”, following that the reserves do not differ in a qualitative way. Note that, the equilibrium strategies of the fossil-fuel owner,  $Q^H = Q_k^1$ , just equal these strategies,  $Q_1^{H**} = Q_1^{op}$ ,  $Q_2^{H**} = Q_2^{op}$ ,  $\tau^{H**} = \tau^{op}$ .

We now describe the driving forces of the symmetric information model.

Neglecting strategic behavior, the fossil-fuel owner, considering the social discount factor, prefers to sell more fossil-fuels in period 1 than in period 2. The social planner, on the other hand, attaches importance not only to the social discount factor but also to the cumulative disutility of carbon emissions. Through the social discount factor the social planner admittedly prefers a higher supply in period 1 than in period 2 alike the fossil-fuel owner. However, the cumulative structure of the environmental damages shifts her preferences. She prefers a lesser supply in period 1 and a higher supply in period 2 compared to the fossil-fuel owner’s preferred supplies.

Now incorporating strategic behavior, the social planner can use the tax as an aid to force the fossil-fuel owner to supply the amounts which are more in her interests, by punishing the fossil-fuel owner with a high tax if he sells too much fossil-fuels in period 1 and rewarding him with a low tax if the supply of period 1 is in her interests. Analyzing the model we show that the regulator’s tax response,  $\tau^{x**}(Q_1)$ , is indeed an increasing function of  $Q_1$ , and equals

$$\tau^{x**}(Q_1) = \max \{ \tau_{op}(Q_1), \tau_b(Q_1, Q^x) \}.^{15} \quad (3.16)$$

The tax, of course, affects the supply of period 2. If the fossil-fuel owner faces a higher tax in period 2, the fossil-fuel price is higher in period 2, and hence he sells a lower amount of his reserves in period 2. In other words, the fossil-fuel owner’s best response,  $Q_2(\tau)$ , is a decreasing function of the tax and it turns out that it equals

---

<sup>15</sup>  $\tau_{op}(Q_1) \equiv 1 + \frac{2dQ_1(1+r) - 2(1-\lambda)r}{3r - 5\lambda r + d(1+r)}$  and  $\tau_b(Q_1, Q^x) \equiv 1 - 2Q^x + 2Q_1$ .

$$Q_2^{x**}(Q_1, \tau) = \min \left\{ \frac{1-\tau}{2}, Q^x - Q_1 \right\}. \quad (3.17)$$

Taking the social planner's tax strategy into account the fossil-fuel owner is going to trade off the effect of the social discount factor on his utility against the effect of the social planner's tax strategy on his utility when choosing his preferred supplies,  $(Q_1^{x**}, Q_2^{x**})$ .

Before we state the results of this section we highlight another two features of our model, which are useful for our analysis.

The first feature is that the regulator profits from higher taxes if the reserves are exhausted in period 2. To understand this consider problem (3.7) under the constraint  $Q_2(\tau) = Q^x - Q_1$ . Given this constraint the regulator's problem becomes

$$\begin{aligned} \tau^{x**}(Q_1) = \operatorname{argmax}_{\tau} \left\{ (1-2\lambda)(Q^x - Q_1) \cdot \tau + \frac{1}{2}(1-\lambda)(Q^x - Q_1)^2 \right. \\ \left. + \lambda(1 - Q^x + Q_1)(Q^x - Q_1) - \frac{1+r}{r} \frac{d}{2}(Q^x)^2 \right\}, \end{aligned} \quad (3.18)$$

which is linear in the tax and has a positive slope. So, if the regulator wants the reserves to be exhausted she profits from setting the highest tax which allows the fossil-fuel owner to exhaust his reserves. In this situation the tax cannot be higher than the described tax as otherwise the fossil-fuel owner's response in period 2 is lower than the remaining reserves and the fossil-fuel reserves will not be exhausted. Solving the above problem (3.18) results in

$$\tau^{x**}(Q_1) = \tau_b(Q_1, Q^x). \quad (3.19)$$

The second feature is that the regulator "discriminates" fossil-fuel owners who have lower reserves left in period 2 against fossil-fuel owners who have higher reserves left in period 2. This is summarized in the next lemma.

**Lemma 3.1** *Given a supply  $0 \leq Q_1 \leq Q^L$  the regulator's tax response always satisfies*

$$\tau^{H**}(Q_1) \leq \tau^{L**}(Q_1).$$

**Proof:** See in appendix C.1. ■

To get an intuition for this result, note that the lesser fossil-fuels are available in period 2 the higher the tax can be to force the fossil-fuel owner to exhaust his reserves. So at least, if the regulator prefers the reserves to be exhausted she sets higher taxes if lesser fossil-fuels are available in period 2. The equilibrium tax function however, as a function of the size of the reserves, does not inherit this property. In equilibrium, fossil-fuel owners with different size of fossil-fuel reserves supply different amounts in period 1 and hence might e.g. have the same amount of fossil-fuels left in period 2.

The following proposition states some features of the symmetric information equilibrium strategies and some comparative statics results with respect to changes in  $d$  and  $\lambda$ .

**Proposition 3.1** *The equilibrium strategies of the low and the high type,  $0 < Q^L < Q^H \leq Q_k^1$ , satisfy*

$$Q_1^{x^{**}} + Q_2^{x^{**}} = Q^x; \quad \tau^{x^{**}} = \tau_b(Q_1^{x^{**}}, Q^x), \quad \text{with } x \in \{L, H\}.$$

Let  $x \in \{L, H\}$ .

1. *Given two different damage coefficients,  $d < \tilde{d}$ , the equilibrium strategies and the threshold,  $Q_k^1$ , relate as follows*

$$Q_k^1(d) \geq Q_k^1(\tilde{d}); \quad Q_1^{x^{**}}(d) \leq Q_1^{x^{**}}(\tilde{d}); \quad Q_2^{x^{**}}(d) \geq Q_2^{x^{**}}(\tilde{d}); \quad \tau^{x^{**}}(d) \leq \tau^{x^{**}}(\tilde{d}).$$

2. *Given two different weights,  $\lambda < \tilde{\lambda}$ , the equilibrium strategies and the threshold,  $Q_k^1$ , relate as follows*

$$Q_k^1(\lambda) \leq Q_k^1(\tilde{\lambda}); \quad Q_1^{x^{**}}(\lambda) \leq Q_1^{x^{**}}(\tilde{\lambda}); \quad Q_2^{x^{**}}(\lambda) \leq Q_2^{x^{**}}(\tilde{\lambda}); \quad \tau^{x^{**}}(\lambda) \geq \tau^{x^{**}}(\tilde{\lambda}).$$

**Proof:** See in appendix C.1. ■

To point 1. If the damage coefficient raises the negative effect of one unit burned fossil-fuels on the environment raises and consequently the dis-utility on the social planner's welfare in period 2 raises, too. This entails a higher tax response by the social planner, which again results in the fossil-fuel owner preferring to supply more in period 1 and less in period 2 compared to his preferred supplies when facing a lower tax.

However, this shift in the fossil-fuel owner's strategy is not really in the interest of the social planner. In contrast, as of the cumulative structure of the environmental costs she prefers a lower supply in period 1 and, if necessary, a higher supply in period 2, when the damage coefficient raises. Yet, she cannot enforce these supplies as she cannot influence the fossil-fuel owner's supply in period 1, directly. But at least –from the social planner's perspective– the total supply decreases with a raise in the damage coefficient as the threshold,  $Q_k^1$ , declines. Consequently less fossil-fuel reserves are going to be exhausted in equilibrium.

To point 2. When  $\lambda$  increases the social planner puts a higher weight on the fossil-fuel owner's utility and a smaller weight on the consumer's utility. In other words, the profit of the fossil-fuel owner has a higher positive effect on the social planner's utility when  $\lambda$  is higher. Hence, if  $\lambda$  increases the social planner's tax response decreases to enable the fossil-fuel owner to make higher profits in period 1 as well as in period 2.

We also verify that if  $\lambda$  raises also the threshold,  $Q_k^1$ , raises, consequently more fossil-fuel reserves are going to be exhausted in equilibrium. Fortunately, this is exactly in the social planner's interests, as with a raise in  $\lambda$  the impact of the fossil-fuel owner's profits compared to the impact of the environmental costs on the social planner's utility gets higher.

### 3.4 The Case of Asymmetric Information

In the general setting the social planner aiming to regulate fossil-fuel consumption as the stock of emissions creates environmental damage does not know how much fossil-fuels,  $Q$ , exactly exist. In contrast to the fossil-fuel owner who knows the level of his reserves,  $Q \in \{Q^L, Q^H\}$ , where  $0 \leq Q^L < Q^H \leq Q_k^1$ . The social planner, however, knows the distribution

$$p_r(Q^L) = \mu, \quad p_r(Q^H) = 1 - \mu, \quad \text{where } 0 < \mu < 1.$$

We define a certain set of monotone belief systems and assume that the allowed belief systems of this game are elements of this set only.

**Definition 3.1** *Let  $\Omega$  denote the set of belief systems consisting of the following elements. Given a probability,  $\mu \in [0, 1]$ , a belief system,  $\mu^*$ , belongs to the set,  $\Omega$ , (i.e.  $\mu^* \in \Omega$ ), if and only if: On the equilibrium path, the beliefs are consistent with Bayes Rule. And off the*

equilibrium path, the beliefs satisfy

$$\begin{aligned}\mu^*(Q^L | Q_1) &= 1 \quad \text{if } Q_1 \leq Q^L, \\ \mu^*(Q^L | Q_1) &= 0 \quad \text{else,}\end{aligned}$$

where  $Q_1 \neq Q_1^{L*}$  and  $Q_1 \neq Q_1^{H*}$ .

We distinguish two kinds of Perfect Bayesian Nash equilibrium strategies. A strategy profile and a belief system  $\mu^* \in \Omega$  which constitute an equilibrium, belong

i) either to the set of separating equilibria: In a separating equilibrium the two types of fossil-fuel owners supply two different amounts of fossil-fuels in period 1. Hence, in a separating equilibrium the social planner can certainly infer the fossil-fuel owner's type after period 1. A separating equilibrium, given the belief system,  $\mu^*$ , is determined by strategies

$$(Q_1^{L*}, Q_2^{L*}, Q_1^{H*}, Q_2^{H*}, \tau^{L*}, \tau^{H*}),$$

with  $Q_1^{L*} \neq Q_1^{H*}$ ;

ii) or to the set of pooling equilibria: In a pooling equilibrium the two types of fossil-fuel owners supply the same amount of fossil-fuels in period 1. So, in a pooling equilibrium the regulator after observing  $Q_1$  does not know for sure which type she is facing. A pooling equilibrium, given the belief system,  $\mu^*$ , is determined by strategies

$$(Q_1^*, Q_2^{L*}, Q_2^{H*}, \tau^*),$$

with  $Q_1^* = Q_1^{L*} = Q_1^{H*}$ .

The remainder of section 3.4 identifies the properties of the best response functions as well as of the equilibrium strategies. Dividing the analysis into two parts, separating and pooling equilibria, each time we start with the fossil-fuel owner's best response in period 2, then we turn to the regulator's best tax response, and finally by using these results we establish the properties of the equilibrium strategies and their conditions to exist. The equilibrium selection follows after that.

### 3.4.1 Separating Equilibria

Let  $\mu^* \in \Omega$ . And let,

$$(Q_1^{L*}, Q_2^{L*}, Q_1^{H*}, Q_2^{H*}, \tau^{L*}, \tau^{H*}), \text{ with } Q_1^{L*} \neq Q_1^{H*},$$

denote a corresponding separating equilibrium strategy profile.

**Best responses in period 2** Given a supply  $Q_1 \leq Q^H$  –on or off the separating equilibrium path– the fossil-fuel owner  $x \in \{L, H\}$  deciding what to supply in period 2,  $Q_2^{x*}(Q_1, \tau)$ , only takes the tax,  $\tau$ , and the size of his remaining reserves,  $Q^x - Q_1$ , into account, like in the symmetric information case. In particular, the belief,  $\mu^*(Q^L | Q_1)$ , does not affect his decisions in period 2. This leads to the conclusion that the fossil-fuel owner’s best response in the asymmetric information case coincides with his best response in the symmetric information case given the right level of reserves. In other words, equation (3.11) equals equation (3.15), and consequently the following lemma is true.

**Lemma 3.2** *The fossil-fuel owner  $x \in \{L, H\}$ ’s best response in period 2,  $Q_2^{x*}(Q_1, \tau)$ , given any  $0 \leq Q_1 \leq Q^x$  and any  $\tau \in [0, 1]$  satisfies*

$$Q_2^{x*}(Q_1, \tau) = Q_2^{x**}(Q_1, \tau). \quad (3.20)$$

**Proof:** The results are straightforward. The proof is therefore omitted. ■

**Best tax response** Given a supply  $Q_1 \leq Q^H$  the regulator deciding which tax,  $\tau^*(Q_1)$ , to set takes the belief,  $\mu^*(Q^L | Q_1)$ , and the corresponding best response,  $Q_2^{x*}(Q_1, \tau^*(Q_1))$ , into account. After observing the supply,  $Q_1$ , –on or off the separating equilibrium path– the regulator believes that she is facing the low type with probability 1 if  $Q_1 \leq Q^L$  and  $Q_1 \neq Q_1^{H*}$ , and she believes that she is facing the high type with probability 1 if  $Q_1 = Q_1^{H*}$  or  $Q^L < Q_1 \leq Q^H$ .

As lemma 3.2 is true, equation (3.10) coincides with equation (3.14) given the right level of reserves. Hence, the regulator responds in the same way as she would do in the symmetric information case being subject to the corresponding level of reserves. More precisely, she responds via  $\tau^{L**}(Q_1)$  (“L-Tax-Response”) if she believes that she is facing the low type, and via  $\tau^{H**}(Q_1)$  (“H-Tax-Response”) if she believes that she is facing the high type. This gives



us the following lemma.

**Lemma 3.3** *The regulator's best tax response,  $\tau^*(Q_1)$ , given any supply,  $0 \leq Q_1 \leq Q^H$ , in period 1 equals*

$$\tau^*(Q_1) = \begin{cases} \tau^{L^{**}}(Q_1) & \text{if } 0 \leq Q_1 \leq Q^L \text{ and } Q_1 \neq Q_1^{H^*} \\ \tau^{H^{**}}(Q_1) & \text{else} \end{cases}. \quad (3.21)$$

**Proof:** The results are straightforward. The proof is therefore omitted. ■

This lemma states that given any separating equilibrium the low type deviating from  $Q_1^{L^*}$  to any strategy  $Q_1 \leq Q^L$ , with  $Q_1 \neq Q_1^{H^*}$ , subsequently faces the L-Tax-Response. And only if he deviates to  $Q_1^{H^*}$  he faces the H-Tax-Response. In contrast to the low type, the high type faces the L-Tax-Response not only for one possible deviation, but for a range of deviations. Given an equilibrium the high type deviating from  $Q_1^{H^*}$  to any strategy  $Q_1 \leq Q^L$  faces the L-Tax-Response and only deviating to any strategies  $Q^L < Q_1 \leq Q^H$ , faces the H-Tax-Response.

**Equilibrium strategies** Before we identify the separating equilibrium strategies let us define some functions to ease the analysis.

The functions,  $U_{H,y}(Q_1)$ ,  $y \in \{L, H\}$ , describe the high type's ex ante utilities at period 1 as functions of  $Q_1$  and of facing either the H-Tax-Response (if  $y = H$ ) or the L-Tax-Response (if  $y = L$ ) and subsequently supplying his corresponding best response in period 2 (see lemma 3.2). The other two functions,  $U_{L,y}(Q_1)$ ,  $y \in \{L, H\}$ , are defined in the same fashion.

Let

$$U_{x,y}(Q_1) \equiv U_F(Q_1, Q_2^{x^{**}}(Q_1, \tau^{y^{**}}(Q_1)), \tau^{y^{**}}(Q_1)), \quad x \in \{L, H\}, \quad y \in \{L, H\}.^{16} \quad (3.22)$$

Maximizing these functions by the first order condition yields  $U_{x,y}(Q_1^{y^{**}}) = \max_{Q_1} U_{x,y}(Q_1)$  for all  $x \in \{L, H\}$  and  $y \in \{L, H\}$ .

Using the definition and applying the previous results (lemma 3.2 and lemma 3.3) we can

---

<sup>16</sup>As the underlying utility function,  $U_F(Q_1, Q_2(Q_1, \tau(Q_1)), \tau(Q_1))$ , is concave in  $Q_1$  the functions,  $U_{x,y}(Q_1)$ ,  $x \in \{L, H\}$ ,  $y \in \{L, H\}$ , are concave in  $Q_1$  as well. So, the inverse functions,  $U_{x,y}^{-1}(\cdot)$ ,  $x \in \{L, H\}$ ,  $y \in \{L, H\}$ , have at most two values.

rewrite the equilibrium conditions: Let  $\mu^* \in \Omega$ . A strategy profile,

$$(Q_1^{L*}, Q_2^{L*}, Q_1^{H*}, Q_2^{H*}, \tau^{L*}, \tau^{H*}), \text{ with } Q_1^{L*} \neq Q_1^{H*},$$

together with  $\mu^*$  constitute a separating equilibrium if and only if for all  $x \in \{L, H\}$ :

$$Q_2^{x*} = Q_2^{x**}(Q_1^{x*}, \tau^{x*}), \quad \tau^{x*} = \tau^{x**}(Q_1^{x*}),$$

and given any deviation  $0 \leq Q_1 \leq Q^x$  in period 1

$$U_{x,x}(Q_1^{x*}) \geq \begin{cases} U_{x,L}(Q_1) & \text{if } 0 \leq Q_1 \leq Q^L \text{ and } Q_1 \neq Q_1^{H*} \\ U_{x,H}(Q_1) & \text{if } Q_1 = Q_1^{H*} \text{ or } Q^L < Q_1 \leq Q^H \end{cases}.^{17} \quad (3.23)$$

The following proposition identifies the separating equilibrium strategies. Following this we give a reasoning for these results.

**Proposition 3.2** *Let  $\mu^* \in \Omega$ .*

1. *If  $Q_1^{H**} > Q^L$  the “symmetric information strategy” profile,*

$$(Q_1^{L**}, Q_2^{L**}, Q_1^{H**}, Q_2^{H**}, \tau^{L**}, \tau^{H**}), \quad (3.24)$$

*and  $\mu^*$  constitute a separating equilibrium.*

2. *If  $Q_1^{H**} \leq Q^L$  the strategy profile,*

$$(Q_1^{L**}, Q_2^{L**}, Q_1^H, Q_2^{H**}(Q_1^H, \tau^{H**}(Q_1^H)), \tau^{L**}, \tau^{H**}(Q_1^H)), \text{ with } Q_1^H \neq Q_1^{L**}, \quad (3.25)$$

---

<sup>17</sup>This condition is a rewritten version of equation (3.8) given the best responses, lemma 3.2 and lemma 3.3. Note that condition (3.23) can be also stated as

$$U_{x,x}(Q_1^{x*}) \geq \max \left\{ \max_{(0 \leq Q_1 \leq Q^L) \vee (Q_1 \neq Q_1^{H*})} U_{x,L}(Q_1); \max_{(Q_1 = Q_1^{H*}) \wedge (Q^L < Q_1 \leq Q^H)} U_{x,H}(Q_1) \right\}.$$

and  $\mu^*$  constitute a separating equilibrium if and only if

$$U_{L,L}(Q_1^{L**}) \geq U_{L,H}(Q_1^H) \quad \text{and} \quad Q_1^H \in [Q_1^{H-}, Q_1^{H+}].^{18} \quad (3.26)$$

There does not exist any other strategy profile which together with  $\mu^*$  constitutes a separating equilibrium.

**Proof:** See in appendix C.2. ■

To understand this proposition note the following.

To the low type. Given any separating equilibrium the low type deviating from  $Q_1^{L*}$  to any strategy off the equilibrium path,  $Q_1 \leq Q^L$ , with  $Q_1 \neq Q_1^{H*}$ , subsequently faces the L-Tax-Response, and only if he deviates to  $Q_1^{H*}$  faces the H-Tax-Response. So, it makes sense that the low type's best supply in period 1 is his symmetric information supply,  $Q_1^{L**}$ , as long as he does not have any incentives to mimic the high type in period 1 by supplying  $Q_1^{H*}$ . If  $Q_1^{H**} > Q^L$  the low type does not have the quantity capacity to mimic the high type. If  $Q_1^{H**} \leq Q^L$  it is the same as assuming,  $U_{L,L}(Q_1^{L**}) \geq U_{L,H}(Q_1^H)$ .

To the high type. In contrast to the low type the high type faces the L-Tax-Response not only for one possible deviation, but for a range of deviations. This makes the analysis of the high type a little bit more complex. However, if  $Q_1^{H**} > Q^L$  the high type always faces the H-Tax-Response after supplying or deviating to  $Q_1^{H**}$ . Recall that the H-Tax-Response is the lowest possible tax response the fossil-fuel owner could face, the L-Tax-Response is higher. So, the high type given the strategies in (3.24) does not have an incentive to deviate in period 1. But given any other strategy profile he always has an incentive to deviate to  $Q_1^{H**}$  in period 1. If, on the other hand,  $Q_1^{H**} \leq Q^L$  condition  $Q_1^H \in [Q_1^{H-}, Q_1^{H+}]$  exactly ensures that the high type does not have an incentive to deviate in period 1.

Analyzing the set of separating equilibria with respect to the probability,  $\mu$ , we find that

---

<sup>18</sup> The thresholds are defined by

$$Q_1^{H-} \equiv \min \left\{ U_{H,H}^{-1} \left( \max[ U_{H,L}(Q_1^{L**}), U_{H,H}(Q^L) ] \right) \right\}$$

and

$$Q_1^{H+} \equiv \max \left\{ U_{H,H}^{-1} \left( \max[ U_{H,L}(Q_1^{L**}), U_{H,H}(Q^L) ] \right) \right\}.$$

Note that the set,  $[Q_1^{H-}, Q_1^{H+}]$ , is not empty as the high type's symmetric information supply,  $Q_1^{H**}$ , always satisfies  $Q_1^{H**} \in [Q_1^{H-}, Q_1^{H+}]$ .

the size of the set of separating equilibria is independent of the probability,  $\mu$ . The following lemma makes the statement preciser.

**Lemma 3.4** *If a strategy profile,  $(Q_1^{L*}, Q_2^{L*}, Q_1^{H*}, Q_2^{H*}, \tau^{L*}, \tau^{H*})$ ,  $Q_1^{L*} \neq Q_1^{H*}$ , and a belief system,  $\mu^* \in \Omega$ , with  $\mu \in [0, 1]$ , constitute a separating equilibrium then the same strategy profile and any other belief system  $\tilde{\mu}^* \in \Omega$ , with  $\tilde{\mu} \in [0, 1]$  and  $\tilde{\mu} \neq \mu$  also constitute a separating equilibrium.*

**Proof:** The result is straightforward. Given a separating strategy profile the regulator's belief after observing any supply  $Q_1$  -on or off the separating equilibrium path- only takes the two values, either  $\mu^*(Q^L | Q_1) = 0$  or  $\mu^*(Q^L | Q_1) = 1$ . So, the separating equilibrium strategies, the deviation incentives and equilibrium conditions do not depend on the probability,  $\mu \in [0, 1]$ . Which means that the probability,  $\mu$ , is, so to say, nonexistent. A formal proof is omitted. ■

### 3.4.2 Pooling Equilibria

Let  $\mu^* \in \Omega$ . And let,

$$(Q_1^*, Q_2^{L*}, Q_2^{H*}, \tau^*),$$

denote a corresponding pooling equilibrium strategy profile.

**Best responses in period 2** Like in the separating equilibrium case, given a supply,  $Q_1 \leq Q^H$  -on or off the pooling equilibrium path- the fossil-fuel owner  $x \in \{L, H\}$  deciding what to supply in period 2 only takes the tax and the size of his remaining reserves into account, like in the symmetric information case. So again, the fossil-fuel owner  $x \in \{L, H\}$ 's best response in period 2 given any  $Q_1 \leq Q^x$  and any  $\tau \in [0, 1]$  satisfies

$$Q_2^{x*}(Q_1, \tau) = Q_2^{x**}(Q_1, \tau). \quad (3.27)$$

**Best tax response** Given a supply  $Q_1 \leq Q^H$  the regulator deciding which tax,  $\tau^*(Q_1)$ , to set takes the belief and the corresponding best response of period 2 into account. After observing the supply,  $Q_1$ , of period 1 -on or off the pooling equilibrium path- the regulator has, in contrast to the separating equilibrium situation, three different beliefs. She believes

that she is facing the low type with probability 1 if  $Q_1 \leq Q^L$  and  $Q_1 \neq Q_1^*$ ; she believes that she is facing the high type with probability 1 if  $Q^L < Q_1 \leq Q^H$ ;<sup>19</sup> and finally she believes that she is facing the low type with probability,  $\mu$ , and the high type with probability,  $1 - \mu$ , if  $Q_1 = Q_1^*$ .

So, the analysis of the regulator's best tax response,  $\tau^*(Q_1)$ , as long as  $Q_1 \neq Q_1^*$  can be simply adopted from the separating equilibrium analysis.

However, if the supply satisfies  $Q_1 = Q_1^*$  the regulator's tax response,  $\tau^*(Q_1)$ , should and does balance out the low and the high type's best responses in period 2,  $Q_2^{x*}(Q_1, \tau)$ , given the probabilities with which they occur,  $\mu^*(Q^x | Q_1)$ . Her tax response lays between her tax response knowing that she is facing the low type and her tax response knowing that she is facing the high type. Let us divide the reasoning into two parts.

We start with a simple case. Assume  $\tau^{L**}(Q_1) = \tau^{H**}(Q_1)$ . The regulator if she would know the fossil-fuel owner's type for sure prefers not to exhaust the reserves independent of the size of the reserves and sets the same tax independent of the type. As the regulator does not distinguish between both types, it is reasonable that also in the asymmetric information case her tax response just equals the symmetric information tax responses,

$$\tau^*(Q_1) = \tau^{L**}(Q_1) = \tau^{H**}(Q_1).$$

Now, assume  $\tau^{L**}(Q_1) \neq \tau^{H**}(Q_1)$ . Note that this can only be the case if at least the low type exhausts his reserves in period 2.<sup>20</sup> We now argue that the pooling tax,  $\tau^*(Q_1)$ , indeed lays in between the two symmetric information tax responses,  $\tau^{H**}(Q_1) \leq \tau^*(Q_1) \leq \tau^{L**}(Q_1)$ .

First,  $\tau^*(Q_1)$  cannot be lower than the H-Tax-Response, as the regulator profits from tax revenues whenever the reserves are exhausted. And if  $\tau^*(Q_1)$  equals or is lower than the H-Tax-Response both types exhaust their reserves. Hence, it is profitable for the regulator to set the highest possible tax,  $\tau^{H**}(Q_1)$ , under these circumstances.

Second,  $\tau^*(Q_1)$  also cannot be higher than the L-Tax-Response, as given any of these taxes both types do not exhaust their reserves, which is not in the regulator's interests, as she prefers at least the low type to exhaust his reserves. So, her tax response definitely forces the low type to exhaust his reserves in period 2.

Now, what is the regulator's trade off whenever  $\tau^{L**}(Q_1) \neq \tau^{H**}(Q_1)$ ? The low type's

---

<sup>19</sup>Note that  $Q_1^* \leq Q^L$  as otherwise the low type cannot supply this amount.

<sup>20</sup>i.e.  $\tau^{L**}(Q_1) \neq \tau^{H**}(Q_1)$  if and only if  $Q_2^{L**}(Q_1, \tau^{L**}(Q_1)) = Q^L - Q_1$ .

strategy in period 2 is independent of the exact tax level as  $\tau \leq \tau^{L**}(Q_1)$ , he always exhaust his reserves. However, as the regulator profits from tax revenues whenever the reserves are exhausted (see equation (3.18)) she prefers to set the highest possible tax,  $\tau^{L**}(Q_1)$ , if she faces the low type. If the regulator instead faces the high type she can effectively influence his supply of period 2 by choosing different tax responses in the given range. Now, recall from the symmetric information case that the regulator facing the high type wants him to supply,  $Q_2^{H**}(Q_1, \tau^{H**}(Q_1))$ , and hence she prefers to set the H-Tax-Response if she faces the high type. So, the regulator setting an optimal tax,  $\tau^*(Q_1)$  –possibly higher than the H-Tax-Response–, trades off, her loss from the high type not supplying the “right” supply, to her gain from the higher tax revenues when the fossil-fuel owner is instead the low type.

**Lemma 3.5** *The regulator’s best tax response,  $\tau^*(Q_1)$ , given any supply,  $0 \leq Q_1 \leq Q^H$ , equals*

$$\tau^*(Q_1) = \begin{cases} \tau^{L**}(Q_1) & \text{if } 0 \leq Q_1 \leq Q^L \text{ and } Q_1 \neq Q_1^* \\ \min \{ \tau^{L**}(Q_1), \max \{ \tau^{H**}(Q_1), \tau_r(Q_1) \} \} & \text{if } Q_1 = Q_1^* \\ \tau^{H**}(Q_1) & \text{if } Q^L < Q_1 \leq Q^H \end{cases} \quad .^{21} \quad (3.28)$$

**Proof:** See in appendix C.2. ■

It follows immediately that for all  $0 \leq Q_1 \leq Q^H$  the tax response,  $\tau^*(Q_1)$ , satisfies

$$\tau^{H**}(Q_1) \leq \tau^*(Q_1) \leq \tau^{L**}(Q_1). \quad (3.29)$$

However, it is not immediately clear if and when the inequalities of equation (3.29) are strict for  $Q_1 = Q_1^*$ . Do parameter constellations exist such that one of these inequalities is an equality? To answer this question let us have a closer look at how  $\tau^*(Q_1^*)$  reacts to changes in the probability,  $\mu$ .

**Lemma 3.6** *1. The tax response,  $\tau^*(Q_1^*)$ , increases with an increase in the probability,  $\mu$ .*

---

<sup>21</sup>  $\tau_r(Q_1) \equiv \tau_{op}(Q_1) + \frac{\mu}{(1-\mu)} \frac{(1-2\lambda)(Q^L - Q_1)}{3r-5\lambda r+d(1+r)}$ .

2. There exist two thresholds,  $0 < \mu_T(Q_1^*) < \mu_{\bar{T}}(Q_1^*) < 1$ , such that

$$\tau^*(Q_1^*) = \begin{cases} \tau^{H^{**}}(Q_1^*) & \text{if and only if } \mu \leq \mu_T(Q_1^*) \\ \tau^{L^{**}}(Q_1^*) & \text{if and only if } \mu \geq \mu_{\bar{T}}(Q_1^*) \end{cases}. \quad (3.30)$$

**Proof:** See in appendix C.2. ■

The intuition of this result is simple. The probability,  $\mu$ , measures how likely the regulator faces the different types. Hence it is reasonable that if the fossil-fuel owner is more likely the high (low) type the regulator ex ante prefers to act as if she is facing the high (low) type more than acting as if she is facing the low (high) type. More precisely, the more likely the regulator faces the low (high) type the closer the tax response,  $\tau^*(Q_1^*)$ , should be to the L-Tax-Response (H-Tax-Response). In other words, the tax response,  $\tau^*(Q_1^*)$ , is an increasing function of the probability,  $\mu$ . In addition, we find that whenever the probability,  $\mu$ , is sufficiently low (high) the regulator does not care at all about the low (high) type's reactions. This result gives us the thresholds.

Let “P-Tax-Response” denote the tax response,  $\tau^*(Q_1)$ , from lemma 3.5 and let “P-Tax” denote the term,  $\min \{ \tau^{L^{**}}(Q_1), \max \{ \tau^{H^{**}}(Q_1), \tau_r(Q_1) \} \}$ .

**Equilibrium strategies** Before we identify the pooling equilibrium strategies let us define two functions which are similar to the functions  $U_{x,y}(Q_1)$ ,  $x \in \{L, H\}$ ,  $y \in \{L, H\}$ .

The function,  $U_{H,*}(Q_1)$ , describes the high type's ex ante utility at period 1 as a function of  $Q_1$  and of facing the P-Tax-Response and subsequently supplying his corresponding best response in period 2. The function,  $U_{L,*}(Q_1)$ , is defined in the same fashion.

Let

$$U_{x,*}(Q_1) \equiv U_F(Q_1, Q_2^{x^{**}}(Q_1, \tau^*(Q_1)), \tau^*(Q_1)), \quad x \in \{L, H\}. \quad (3.31)$$

Using this definition and the previous results (equation 3.27 and lemma 3.5) we can again rewrite the equilibrium conditions: Let  $\mu^* \in \Omega$ . A strategy profile,

$$(Q_1^*, Q_2^{L^*}, Q_2^{H^*}, \tau^*),^{22}$$

---

<sup>22</sup>Recall that a necessary condition is  $Q_1^* \leq Q^L$ .

together with  $\mu^*$  constitute a pooling equilibrium if and only if for all  $x \in \{L, H\}$ :

$$Q_2^{x*} = Q_2^{x**}(Q_1^*, \tau^*), \quad \tau^* = \tau^*(Q_1^*),$$

and given any deviation  $0 \leq Q_1 \leq Q^x$  in period 1

$$U_{x,*}(Q_1^*) \geq \begin{cases} U_{x,L}(Q_1) & \text{if } 0 \leq Q_1 \leq Q^L \text{ and } Q_1 \neq Q_1^* \\ U_{x,H}(Q_1) & \text{if } Q^L < Q_1 \leq Q^H \end{cases}. \quad (3.32)$$

The following proposition states the pooling equilibrium strategies. Following this we give a reasoning for these results.

**Proposition 3.3** *Let  $\mu^* \in \Omega$ .*

1. *If  $Q_1^{H**} > Q^L$  there does not exist a strategy profile which together with  $\mu^*$  constitutes a pooling equilibrium.*
2. *If  $Q_1^{H**} \leq Q^L$  the following strategy profile,<sup>24</sup>*

$$(Q_1^*, Q_2^{L**}(Q_1^*, \tau^*(Q_1^*)), Q_2^{H**}(Q_1^*, \tau^*(Q_1^*)), \tau^*(Q_1^*)), \quad (3.33)$$

*and  $\mu^*$  constitute a pooling equilibrium if and only if*

$$U_{L,*}(Q_1^*) \geq U_{L,L}(Q_1^{L**}) \text{ and } U_{H,*}(Q_1^*) \geq \max\{U_{H,L}(Q_1^{L**}), U_{H,H}(Q^L)\}. \quad (3.34)$$

*There does not exist any other strategy profile which together with  $\mu^*$  constitutes a pooling equilibrium.*

**Proof:** See in appendix C.2. ■

To point 1. As long as the high type's symmetric information supply satisfies  $Q_1^{H**} > Q^L$ , no pooling equilibrium exists. The reason is simple. If  $Q_1^{H**} > Q^L$  the high type's best

<sup>23</sup>Condition (3.32) can also be stated as

$$U_{x,*}(Q_1^*) \geq \max \left\{ \max_{0 \leq Q_1 \leq Q^L} U_{x,L}(Q_1); \max_{Q^L < Q_1 \leq Q^H} U_{x,H}(Q_1) \right\}.$$

<sup>24</sup>In the appendix we show that the sets of low and high types with these properties are not empty.



supply in period 1 is  $Q_1^{H**}$ . Indeed, supplying any other amount in period 1 he always has an incentive to deviate to  $Q_1^{H**}$  as, i) he always faces the H-Tax-Response after deviating to  $Q_1^{H**}$ , ii)  $U_{H,H}(Q_1^{H**}) = \max_{Q_1} U_{H,H}(Q_1)$ , and iii) the H-Tax-Response is always smaller than the L-Tax-Response, i.e.  $U_{H,L}(Q_1) \leq U_{H,H}(Q_1)$  for all  $Q_1 \leq Q^H$ . The low type, however, does not have the quantity capacity to mimic the high type and therefore no pooling equilibrium exists.

To point 2. Recall that the low type deviating from  $Q_1^*$  always faces the L-Tax-Response. So, a sufficient condition for the low type not having an incentive to deviate to any off equilibrium supply is that he does not have an incentive to deviate to his symmetric information supply of period 1 as  $U_{L,L}(Q_1^{L**}) = \max_{Q_1} U_{L,L}(Q_1)$ . So, condition (3.32) for  $x = L$  is simply equivalent to condition  $U_{L,*}(Q_1^*) \geq U_{L,L}(Q_1^{L**})$ .

To the high type. A sufficient condition for the high type not having an incentive to deviate is that he does not have an incentive to deviate from  $Q_1^*$  to  $Q^L$  or to  $Q_1^{L**}$ . Indeed, as the high type's utility function  $U_{H,H}(Q_1)$  is concave in  $Q_1$  and is maximized at  $Q_1^{H**} \leq Q^L$  the high type's highest utility off the equilibrium path facing the H-Tax-Response is on the boundary,  $U_{H,H}(Q^L)$ . And, his highest utility of facing the the L-Tax-Response is  $U_{H,L}(Q_1^{L**}) = \max_{Q_1} U_{H,L}(Q_1)$ .

Comparing the high type's condition in the separating condition (3.26) to his condition in the pooling condition (3.34) and as  $U_{H,*}(Q_1^*) \leq U_{H,H}(Q_1^*)$  it follows that if  $Q_1^*$  satisfies the high type's condition in (3.34) it also satisfies his condition in (3.26). So, it follows that  $Q_1^* \in [Q_1^{H-}, Q_1^{H+}]$  is a necessary condition for a strategy profile to constitute any kind of equilibrium –pooling or separating.

In order to analyze the set of pooling equilibria in more detail let us have a look at the set with respect to the probability,  $\mu$ . We find that the size of the set of pooling equilibria shrinks with a decrease in the probability,  $\mu$ , whenever  $Q_1^{H**} \leq Q^L$ . The following lemma makes the statement preciser for these cases,  $Q_1^{H**} \leq Q^L$ .

**Lemma 3.7** *If a strategy profile,  $(Q_1^*, Q_2^{L*}, Q_2^{H*}, \tau^*)$ , and a belief system,  $\mu^* \in \Omega$ , with  $\mu \in [0, 1]$ , constitute a pooling equilibrium then the same strategy profile and any other belief system,  $\tilde{\mu}^* \in \Omega$ , with  $0 < \tilde{\mu} \leq \mu$  also constitute a pooling equilibrium. However, this statement is not true for belief systems,  $\tilde{\mu}^* \in \Omega$ , with  $\tilde{\mu} > \mu$ .*

**Proof:** See in appendix C.2. ■

In contrast to a separating equilibrium the probability,  $\mu$ , plays an existent role in a pooling equilibrium. Indeed, after observing the supply,  $Q_1^*$ , the regulator's belief is  $\mu^*(Q^L | Q_1^*) = \mu$ . Now, assume that a strategy profile,  $(Q_1^*, Q_2^{L*}, Q_2^{H*}, \tau^*)$ , and a belief system,  $\mu^* \in \Omega$ , with  $\mu \in [0, 1]$  constitute a pooling equilibrium. As the tax response,  $\tau^*(Q_1^*)$ , increases in  $\mu$  (see lemma 3.6) the fossil-fuel owner's deviation incentives from  $Q_1^*$  to any other supply  $Q_1 \leq Q^H$  increase (decrease) with an increase (decrease) in  $\mu$ . Or, in other words, given a supply,  $Q_1^*$ , condition (3.34) is harder (easier) to met when  $\mu$  increases (decreases). To understand this note that with an increase (decrease) in  $\mu$  the function,  $U_{x,*}(Q_1^*)$ , decreases (increases) but the two functions,  $U_{H,L}(Q_1^{L**})$  and  $U_{H,H}(Q^L)$ , stay constant. This in turn implies that the size of the set of pooling equilibria shrinks with an increase in  $\mu$ .

### 3.4.3 Equilibrium Selection

In order to select from this vast number of equilibria we apply the Cho-Kreps' Intuitive Criterion.

Whenever  $Q_1^{H**} > Q^L$ , the only equilibrium which exists is the separating equilibrium given the strategies of the symmetric information equilibria (see proposition 3.2). Hence, for the remainder of this section we only deal with the cases,  $Q_1^{H**} \leq Q^L$ .

**Lemma 3.8** *Let  $\mu^* \in \Omega$ .*

1. *Assume  $\mu \leq \mu_T(Q_1^{H**})$ . The strategy profile,*

$$(Q_1^{H**}, Q^L - Q_1^{H**}, Q_2^{H**}, \tau^{H**}),^{25} \quad (3.35)$$

*and  $\mu^*$  constitute a pooling equilibrium. This equilibrium is the only equilibrium which survives the Cho-Kreps' Intuitive Criterion.*

2. *Assume  $\mu > \mu_T(Q_1^{H**})$ . There does not exist any equilibrium which survives the Cho-Kreps' Intuitive Criterion.*

**Proof:** See in appendix C.2. ■

---

<sup>25</sup>Note that  $Q_2^{L**}(Q_1^{H**}, \tau^{H**}) = Q^L - Q_1^{H**}$ . This is true as even the high type, owning more fossil-fuel reserves, exhausts his reserves,  $Q_2^{H**} = Q^H - Q_1^{H**}$ , given the supply  $Q_1^{H**}$  of period 1 and the tax  $\tau^{H**}$ .

Before we give a reasoning for this lemma note the following. If  $\mu \leq \mu_T(Q_1^{H^{**}})$  the separating condition (3.26) for  $Q_1^H = Q_1^{H^{**}}$  and the pooling condition (3.34) for  $Q_1^* = Q_1^{H^{**}}$  are complements as  $U_{x,*}(Q_1^{H^{**}}) = U_{x,H}(Q_1^{H^{**}})$  for all  $x \in \{L, H\}$ . Therefore, the separating strategy profile  $(Q_1^{L^{**}}, Q_2^{L^{**}}, Q_1^{H^{**}}, Q_2^{H^{**}}, \tau^{L^{**}}, \tau^{H^{**}})$  together with  $\mu^*$  does not constitute an equilibrium.

To point 1 of the lemma. We first show that the stated strategies together with  $\mu^*$  indeed constitute a pooling equilibrium as long as  $\mu \leq \mu_T(Q_1^{H^{**}})$ . Recall that  $U_{L,H}(Q_1^{H^{**}}) = \max_{Q_1} U_{L,H}(Q_1)$ , i.e. in particular that  $U_{L,H}(Q_1^{L^{**}}) \leq U_{L,H}(Q_1^{H^{**}})$ . Additionally, as the H-Tax-Response is smaller than the L-Tax-Response the following inequality is also true,  $U_{L,L}(Q_1^{L^{**}}) \leq U_{L,H}(Q_1^{L^{**}})$ . Both results together give us the sufficient condition (3.34) for the above stated strategy profile and  $\mu^*$  to constitute a pooling equilibrium.

This pooling equilibrium survives the Cho-Kreps' Intuitive Criterion because no type can strategically set himself better off by supplying an off equilibrium supply,  $Q_1 \neq Q_1^{H^{**}}$ , even if he would consequently always face the lowest possible tax response, which is the H-Tax-Response.<sup>26</sup> Indeed, this is true as  $U_{x,*}(Q_1^{H^{**}}) = U_{x,H}(Q_1^{H^{**}}) \geq U_{x,H}(Q_1)$  for all  $Q_1 \leq Q_1^H$  and  $x \in \{L, H\}$ .

We now deal with all other existing equilibria whether of point 1 or of point 2. We argue that all these equilibria can be eliminated by the Cho-Kreps' Intuitive Criterion. In the appendix we show that given an equilibrium –pooling or separating– there exist an off equilibrium path amount  $Q_1 \leq Q_1^H$  such that, i) the regulator's reasonable belief after observing this amount is to believe that she is facing the high type,  $\mu^*(Q^L | Q_1) = 0$ , and ii) the high type's ex ante utility at period 1 is strictly higher given this amount and the corresponding tax response,  $\tau^{H^{**}}(Q_1)$ , than given the equilibrium strategies. It follows that the equilibrium does not survive the Cho-Kreps' Intuitive Criterion.

We show these statements by pointing out that there exist an off equilibrium path supply  $Q_1 \leq Q_1^H$  such that, i) the low type does not have an incentive to deviate from his equilibrium supply in period 1 to this specific supply,  $Q_1$ , even if he would consequently face the H-Tax-Response, and ii) the high type instead would have an incentive to deviate from his equilibrium supply of period 1 to this specific supply,  $Q_1$ , if he would consequently face the

---

<sup>26</sup>This tax response is the consequence of an off equilibrium path belief which is in the fossil-fuel owner's "best interest", namely the regulator believes that she is facing the high type. So, even if the reasonable off equilibrium path belief is exactly as described the fossil-fuel owner cannot set himself better off by deviating from  $Q_1^{H^{**}}$ .

H-Tax-Response. These two results indeed imply the above two statements.

Taking the results of lemma 3.8 into account we assume for the remainder of this chapter that  $\mu \leq \mu_T(Q_1^{H**})$ .

### 3.5 Welfare Analysis

The welfare analysis is simple whenever  $Q_1^{H**} > Q^L$ . The strategies of the low and the high type's symmetric information equilibria coincide with the strategies of the asymmetric information equilibrium strategy profile given any belief system,  $\mu^* \in \Omega$ , (see proposition 3.2). Hence, the regulator and both types of fossil-fuel owners are equally well off in the symmetric information case as well as in the asymmetric information case. Therefore, from now on we only deal with the cases,  $Q_1^{H**} \leq Q^L$ . Recall that  $\mu \leq \mu_T(Q_1^{H**})$ .

The following proposition states that both types of fossil-fuel owners prefer the pooling equilibrium from lemma 3.8 to all other existing equilibria.

**Proposition 3.4** *Let  $\mu^* \in \Omega$ . Given any equilibrium,  $(Q_1^{L*}, Q_2^{L*}, Q_1^{H*}, Q_2^{H*}, \tau^{L*}, \tau^{H*})$  and  $\mu^*$ , –pooling or separating– the following two inequalities always hold,<sup>27</sup>*

$$U_{L,H}(Q_1^{H**}) \geq U_F(Q_1^{L*}, Q_2^{L*}, \tau^{L*}) \quad (3.36)$$

and

$$U_{H,H}(Q_1^{H**}) \geq U_F(Q_1^{H*}, Q_2^{H*}, \tau^{H*}). \quad (3.37)$$

**Proof:** See in appendix C.3. ■

Note that both types of fossil-fuel owners also prefer the pooling equilibrium from lemma 3.8 to their own symmetric information equilibrium. Indeed, the high type's ex ante utilities at period 1 given both information structures are equal as his strategies coincide in both equilibria. The low type's ex ante utility at period 1 is higher given the asymmetric information equilibrium from lemma 3.8 than given his symmetric information equilibrium as he prefers to mimic the high type when he has the chance to and  $U_{L,L}(Q_1^{L**}) \leq U_{L,H}(Q_1^{H**})$  (see proof of lemma 3.8).

<sup>27</sup>Note that if  $\mu \leq \mu_T(Q_1^{H**})$  the following equality is true  $U_{x,H}(Q_1^{H**}) = U_{x,*}(Q_1^{H**})$ , for all  $x \in \{L, H\}$ .

The reasoning of this proposition is as follows. To the high type. The high type's analysis is simple. In any equilibrium –pooling or separating– the high type's equilibrium tax is either the H-Tax-Response or the P-Tax-Response which is higher than the later. So,  $U_{H,*}(Q_1) \leq U_{H,H}(Q_1)$ , for all  $Q_1 \leq Q^H$ . Together with the analysis of the symmetric information case,  $U_{H,H}(Q_1) \leq U_{H,H}(Q_1^{H**})$ , for all  $0 \leq Q_1 \leq Q^H$ , the result follows immediately.

To the low type. Assume there exist a strategy profile,  $(Q_1^{L*}, Q_2^{L*}, Q_1^{H*}, Q_2^{H*}, \tau^{L*}, \tau^{H*})$ , which together with  $\mu^*$  constitutes a separating equilibrium. In any separating equilibrium the low type's strategies are just his symmetric information equilibrium strategies and hence his ex ante utility at period 1 in a separating equilibrium is  $U_{L,L}(Q_1^{L**})$ . As the strategy profile given in lemma 3.8 together with  $\mu^*$  constitutes a pooling equilibrium the following inequality is true,  $U_{L,L}(Q_1^{L**}) \leq U_{L,H}(Q_1^{H**})$ , i.e. the low type prefers to mimic the high type. This is the same as saying that the low type prefers the pooling equilibrium from lemma 3.8 to any other existing separating equilibrium.

Assume there exist another strategy profile,  $(Q_1^*, Q_2^{L*}, Q_2^{H*}, \tau^*)$ , with  $Q_1^* \neq Q_1^{H**}$ , which together with  $\mu^*$  constitutes a pooling equilibrium. As the P-Tax-Response is higher than the H-Tax-Response it follows that  $U_{L,*}(Q_1^*) \leq U_{L,H}(Q_1^*)$ . And as the high type's symmetric information supply maximizes the low type's utility given the H-Tax-Response,  $U_{L,H}(Q_1^{H**}) = \max_{Q_1} U_{L,H}(Q_1)$ , the result follows immediately. And hence the statement of the proposition is also true for the low type.

**The regulator's welfare** We compare the regulator's ex ante expected utilities at period 0 under both information structures. Recall that  $Q_1^{H**} \leq Q^L$  and  $\mu \leq \mu_T(Q_1^{H**})$ .

The regulator's ex ante expected utility at period 0 in the symmetric information case is simply

$$\mathbb{E}U_{sym}(Q^L, Q^H) \equiv \mu \cdot U_S(Q_1^{L**}, Q_2^{L**}, \tau^{L**}) + (1 - \mu) \cdot U_S(Q_1^{H**}, Q_2^{H**}, \tau^{H**}).^{28} \quad (3.38)$$

Her ex ante expected utility at period 0 given the pooling equilibrium in lemma 3.8 is

$$\mathbb{E}U_{asy}(Q^L, Q^H) \equiv \mu \cdot U_S(Q_1^{H**}, Q^L - Q_1^{H**}, \tau^{H**}) + (1 - \mu) \cdot U_S(Q_1^{H**}, Q_2^{H**}, \tau^{H**}) \quad (3.39)$$

Note that only the low type's strategies differ.

---

<sup>28</sup>Recall that we assume that the regulator faces the fossil-fuel owner  $Q^L$  with probability  $\mu$  and the fossil-fuel owner  $Q^H$  with probability  $1 - \mu$ .

Our aim is to show that the regulator, as it sometimes appears in signaling games, can prefer an asymmetric information equilibrium to the symmetric information cases. The rough intuition is simple. In the symmetric information case the regulator is not able to implement her preferred supplies of both periods as she is only able to influence the supply of period 2 directly through the tax,  $\tau$ . Therefore there exist supply compositions which are more profitable for the regulator. If one of these compositions together with the corresponding tax responses “constitute” an asymmetric information equilibrium the regulator indeed prefers the latter equilibrium to the symmetric information situations.

Analyzing the regulator’s welfare we find that also she can ex ante prefer the asymmetric information equilibrium from lemma 3.8 to the symmetric information equilibria.

**Proposition 3.5** *Let  $\mu^* \in \Omega$ . There always exist pairs  $Q^L$  and  $Q^H$ , with  $Q^L < Q^H$ , such that*

$$\mathbb{E}U_{asy}(Q^L, Q^H) > \mathbb{E}U_{sym}(Q^L, Q^H). \quad (3.40)$$

**Proof:** See in appendix C.3. ■

### 3.6 Conclusion

We have analyzed Pigovian taxation in a natural resources setup with asymmetric information on the size of the reserves. Our main result shows that not only the informed owner, but also the controlling regulator, can profit from her lack of information.

First, the symmetric information analysis reveals that if the low and the high type supply the same amount in the first period, the regulator’s tax response facing the low type is higher than her tax response facing the high type. We then show that the regulator’s tax response, whenever she does not know which type she is facing, is a weighted average of the symmetric information tax responses. Consequently, only the low type might have incentives to mimic the other type. The low type successfully pools with the high type in equilibrium if, in addition, he has the capacity to mimic the high type in the first period. Our results, therefore, suggest that a threat of a future carbon taxation can support over-reporting as claimed by empirical studies.

The social planner can profit from asymmetric information for the following reasons. In our model the social planner is only able to set a tax in the second period, hence cannot

influence the supply of the first period directly, and therefore is not able to implement her preferred supplies of both periods in the symmetric information case. As a consequence, there exist supply compositions which are preferred by the regulator to the symmetric information equilibrium supplies. If these compositions together with the corresponding remaining strategies constitute an asymmetric information equilibrium the regulator indeed prefers this equilibrium to the corresponding symmetric information equilibria.





# Appendix A

## Managerial Incentive Problems and Return Distributions

### A.1 The Problem of Pure Moral Hazard

To ease notation, we introduce the following variables, that are functions of the underlying parameters.

**Definition A.1**

$$\begin{aligned}\eta &\equiv \frac{2(1-\delta)}{2(1-\lambda)-\delta}; \Delta \equiv \lambda^{\frac{2\lambda}{2(1-\lambda)-\delta}} \delta^{\frac{\delta}{2(1-\lambda)-\delta}}; \\ \Gamma &\equiv \left(\frac{1}{2} + \frac{\lambda}{\delta}\right) \left(\delta^{1-\lambda} \lambda^\lambda\right)^{\frac{2}{2(1-\lambda)-\delta}}; \Lambda \equiv \left(\delta^{1-\lambda} \lambda^\lambda\right)^{\frac{1}{2(1-\lambda)-\delta}}.\end{aligned}$$

Moreover,

$$\theta \equiv r \cdot s,$$

where

$$r \equiv a^{\frac{-\delta}{2(1-\lambda)-\delta}} \text{ and } s \equiv c^{\frac{-2\lambda}{2(1-\lambda)-\delta}}.$$

**Proof of Proposition 1.1:** We first establish part ii, then parts i), iii) and iv).

ii) It is straightforward to compute the interior solution from the first-order conditions.

For a given share  $\alpha$ , the agent's effort and volatility choice are given by

$$\sigma(\alpha, \mathbf{t}) = \Lambda r^{\frac{1-\lambda}{\delta}} s^{\frac{1}{2}} \alpha^{\frac{2\lambda-1}{2(1-\lambda)-\delta}} \quad (\text{A.1})$$

and

$$e(\alpha, \mathbf{t}) = \frac{\lambda}{\delta} r^{\frac{2(1-\lambda)-\delta}{-\delta}} s^{\frac{2(1-\lambda)-\delta}{2\lambda}} \alpha^2 \sigma(\alpha, \mathbf{t})^2. \quad (\text{A.2})$$

As shown in the main text, the optimal share  $\alpha^*$  is given by  $\alpha^* = \underline{\alpha} \equiv \frac{\eta-1}{\eta} \frac{\Delta}{\Gamma}$ .

$\sigma(\alpha, \mathbf{t})$  defined by (A.1) is decreasing in  $a$  and  $c$ . The implied mean is

$$\mu(e(\underline{\alpha}, \mathbf{t}), \sigma(\underline{\alpha}, \mathbf{t})) = \left(\frac{\lambda}{\delta}\right)^\lambda \Lambda^{2\lambda+\delta} \underline{\alpha}^{\frac{2\lambda-\delta}{2(1-\lambda)-\delta}} r s,$$

a decreasing function of  $a$  and  $c$ .

i) The function  $\Phi(a)$  is defined by the condition  $\sigma(\underline{\alpha}, \mathbf{t}) = \bar{\sigma}$ , that is

$$\Lambda r^{\frac{1-\lambda}{\delta}} s^{\frac{1}{2}} \left(\frac{2\lambda-\delta}{2(1-\delta)} \frac{\Delta}{\Gamma}\right)^{\frac{2\lambda-1}{2(1-\lambda)-\delta}} = \bar{\sigma}. \quad (\text{A.3})$$

Solving for  $c$ , we obtain

$$c = \Phi(a) \equiv \Lambda^{\frac{2(1-\lambda)-\delta}{\lambda}} \left(\frac{2\lambda-\delta}{2(1-\delta)} \frac{\Delta}{\Gamma}\right)^{\frac{2\lambda-1}{\lambda}} a^{\frac{\lambda-1}{\lambda}} \bar{\sigma}^{-\frac{2(1-\lambda)-\delta}{\lambda}}.$$

iii) For values of  $a$  and  $c$  such that  $c < \Phi(a)$ , it is not optimal to implement  $\sigma(\alpha, \mathbf{t})$  defined by (1.3). Instead, the relevant implementation constraint becomes (1.4). Two possibilities arise. Firstly, (1.4) holds as an equality - which is the case for values of  $\mathbf{t}$  close to the locus defined by  $c = \Phi(a)$ . Since effort is always interior, we can compute  $\alpha^*(\mathbf{t})$  from (1.4) as an equality:

$$\left(\frac{\lambda}{\delta} r^{\frac{2(1-\lambda)-\delta}{-\delta}} s^{\frac{2(1-\lambda)-\delta}{2\lambda}} \alpha(\mathbf{t})^2 \bar{\sigma}^2\right)^\lambda \delta \bar{\sigma}^{\delta-1} - a \alpha(\mathbf{t}) \bar{\sigma} = 0,$$

which yields

$$\alpha^*(\mathbf{t}) = \left(\lambda^\lambda \delta^{1-\lambda} a^{\lambda-1} c^{-\lambda} \bar{\sigma}^{\delta-2(1-\lambda)}\right)^{\frac{1}{1-2\lambda}}, \quad (\text{A.4})$$

a decreasing function of  $a$  and  $c$ . Clearly also, when substituting for  $c = \Phi(a)$  into (A.4), we obtain  $\alpha^*(\mathbf{t}) = \underline{\alpha}$ . Thus, the solution is continuous at the boundary  $c = \Phi(a)$  separating the parameter values that give rise to interior and corner solutions, respectively. The implied mean is

$$\mu(e(\alpha^*(\mathbf{t}), \mathbf{t}), \bar{\sigma}) = \lambda^{\frac{\lambda}{1-2\lambda}} \delta^{\frac{\lambda}{1-2\lambda}} \bar{\sigma}^{\frac{\delta-2\lambda}{1-2\lambda}} a^{-\frac{\lambda}{1-2\lambda}} c^{-\frac{\lambda}{1-2\lambda}},$$

a decreasing function of  $a$  and  $c$ .

Secondly, it can be the case that (1.4) is satisfied automatically - i.e. holds as a strict inequality. Solving the first-order condition for effort for the optimal level of effort for any given contract  $\alpha$ , we obtain

$$e(\alpha, \mathbf{t}) = \lambda^{\frac{1}{1-\lambda}} c^{\frac{-1}{1-\lambda}} \bar{\sigma}^{\frac{\delta}{1-\lambda}} \alpha^{\frac{1}{1-\lambda}}.$$

This implies that the equilibrium mean of the return distribution is

$$\mu^* = \mu(e(\alpha, \mathbf{t}), \bar{\sigma}) = \lambda^{\frac{\lambda}{1-\lambda}} c^{\frac{-1}{1-\lambda}} \bar{\sigma}^{\frac{\delta}{1-\lambda}} \alpha^{\frac{\lambda}{1-\lambda}}.$$

Hence, the principal's problem becomes

$$\max_{\alpha} \left\{ \lambda^{\frac{\lambda}{1-\lambda}} c^{\frac{-\lambda}{1-\lambda}} \bar{\sigma}^{\frac{\delta}{1-\lambda}} \alpha^{\frac{\lambda}{1-\lambda}} - \frac{a}{2} \bar{\sigma}^2 \alpha^2 - \lambda^{\frac{1}{1-\lambda}} c^{\frac{-\lambda}{1-\lambda}} \bar{\sigma}^{\frac{\delta}{1-\lambda}} \alpha^{\frac{1}{1-\lambda}} \right\} \quad (\text{A.5})$$

The solution satisfies the first-order condition

$$\lambda^{\frac{1}{1-\lambda}} \frac{1}{1-\lambda} c^{\frac{-\lambda}{1-\lambda}} \bar{\sigma}^{\frac{\delta}{1-\lambda}} \alpha^* (\mathbf{t})^{\frac{2\lambda-1}{1-\lambda}} - a \bar{\sigma}^2 \alpha^* (\mathbf{t}) - \lambda^{\frac{1}{1-\lambda}} \frac{1}{1-\lambda} c^{\frac{-\lambda}{1-\lambda}} \bar{\sigma}^{\frac{\delta}{1-\lambda}} \alpha^* (\mathbf{t})^{\frac{\lambda}{1-\lambda}} = 0. \quad (\text{A.6})$$

By the second-order condition and the fact that the left-hand side of (A.6) is decreasing in  $a$  and  $c$ , we note that the optimal  $\alpha^* (\mathbf{t})$  is decreasing in  $a$  and  $c$ . Moreover,  $\sigma^* = \bar{\sigma}$  over this range. Note again that  $\mu(e(\alpha^* (\mathbf{t}), \mathbf{t}), \bar{\sigma})$  is decreasing in  $a$  and  $c$ .

Finally, suppose that (1.4) holds as an equality at  $\hat{\mathbf{t}}$  and as a strict inequality at  $\tilde{\mathbf{t}}$  and choose  $\hat{\mathbf{t}}$  and  $\tilde{\mathbf{t}}$  arbitrarily close to each other. Then it must be the case that  $\alpha^* (\tilde{\mathbf{t}}) \geq \alpha^* (\hat{\mathbf{t}})$ . To see this, suppose we had  $\alpha^* (\tilde{\mathbf{t}}) < \alpha^* (\hat{\mathbf{t}})$ . However, then the agent would have a strictly higher incentive to choose marginal incentive to increase  $\sigma$  at  $\hat{\mathbf{t}}$ , contradicting that (1.4) holds as an equality at  $\hat{\mathbf{t}}$  and as a strict inequality at  $\tilde{\mathbf{t}}$ .

iv) This follows from the continuity of the solution at the boundary separating the two sets  $\mathbf{T}_C$  and  $\mathbf{T}_I$  in conjunction with the comparative statics properties of the optimal share  $\alpha^* (\mathbf{t})$  for  $\mathbf{t} \in \mathbf{T}_C$ . ■

**Proof of Proposition 1.2:** Let  $P(\mathbf{T}_C) \equiv \Pr(\mathbf{t} \in \mathbf{T}_C)$  and let  $P(\mathbf{T}_I) \equiv \Pr(\mathbf{t} \in \mathbf{T}_I)$ .

$$COV(\tilde{\alpha}^*, \tilde{\sigma}^*) = \mathbb{E}[\tilde{\alpha}^* \tilde{\sigma}^*] - \mathbb{E}[\tilde{\alpha}^*] \mathbb{E}[\tilde{\sigma}^*].$$

We can rewrite this as

$$\begin{aligned} COV(\tilde{\alpha}^*, \tilde{\sigma}^*) &= \mathbb{E}[\tilde{\alpha}^* \tilde{\sigma}^* | \mathbf{T}_C] P(\mathbf{T}_C) + \mathbb{E}[\tilde{\alpha}^* \tilde{\sigma}^* | \mathbf{T}_I] P(\mathbf{T}_I) \\ &\quad - (\mathbb{E}[\tilde{\alpha}^* | \mathbf{T}_C] P(\mathbf{T}_C) + \mathbb{E}[\tilde{\alpha}^* | \mathbf{T}_I] P(\mathbf{T}_I)) \cdot \\ &\quad (\mathbb{E}[\tilde{\sigma}^* | \mathbf{T}_C] P(\mathbf{T}_C) + \mathbb{E}[\tilde{\sigma}^* | \mathbf{T}_I] P(\mathbf{T}_I)). \end{aligned}$$

Simplifying, we obtain

$$\begin{aligned} COV(\tilde{\alpha}^*, \tilde{\sigma}^*) &= \bar{\sigma} \mathbb{E}[\tilde{\alpha}^* | \mathbf{T}_C] P(\mathbf{T}_C) + \underline{\alpha} \mathbb{E}[\tilde{\sigma}^* | \mathbf{T}_I] P(\mathbf{T}_I) \\ &\quad - (\mathbb{E}[\tilde{\alpha}^* | \mathbf{T}_C] P(\mathbf{T}_C) + \underline{\alpha} P(\mathbf{T}_I)) (\bar{\sigma} P(\mathbf{T}_C) + \mathbb{E}[\tilde{\sigma}^* | \mathbf{T}_I] P(\mathbf{T}_I)). \end{aligned}$$

Multiplying out and rearranging yields

$$\begin{aligned} COV(\tilde{\alpha}^*, \tilde{\sigma}^*) &= \bar{\sigma} \mathbb{E}[\tilde{\alpha}^* | \mathbf{T}_C] P(\mathbf{T}_C) (1 - P(\mathbf{T}_C)) + \underline{\alpha} \mathbb{E}[\tilde{\sigma}^* | \mathbf{T}_I] P(\mathbf{T}_I) (1 - P(\mathbf{T}_I)) \\ &\quad - \underline{\alpha} \bar{\sigma} P(\mathbf{T}_I) P(\mathbf{T}_C) - \mathbb{E}[\tilde{\sigma}^* | \mathbf{T}_I] \mathbb{E}[\tilde{\alpha}^* | \mathbf{T}_C] P(\mathbf{T}_I) P(\mathbf{T}_C). \end{aligned}$$

Using the fact that  $P(\mathbf{T}_I) = 1 - P(\mathbf{T}_C)$  we can rewrite this into

$$COV(\tilde{\alpha}^*, \tilde{\sigma}^*) = (\bar{\sigma} - \mathbb{E}[\tilde{\sigma}^* | \mathbf{T}_I]) P(\mathbf{T}_I) (\mathbb{E}[\tilde{\alpha}^* | \mathbf{T}_C] - \underline{\alpha}) P(\mathbf{T}_C) > 0,$$

where the conclusion follows from the facts that  $\bar{\sigma} > \mathbb{E}[\tilde{\sigma}^* | \mathbf{T}_I]$  and  $\mathbb{E}[\tilde{\alpha}^* | \mathbf{T}_C] > \underline{\alpha}$ .

Part ii) is a trivial consequence of the fact that the functions  $\alpha^*(\mathbf{t})$ ,  $\mu^*(\mathbf{t})$ , and  $\sigma^*(\mathbf{t})$  are comonotone.

Part iii): Since  $\tilde{\mu}^*$  is strictly decreasing in each variable  $a$  and  $c$ , fixing one of these variables, one can write  $\tilde{\alpha}^*$  as a (weakly increasing) function of  $\tilde{\mu}^*$ . The function  $\tilde{\alpha}^*(\tilde{\mu}^*)$  is nondecreasing in  $\tilde{\mu}^*$ , because both  $\tilde{\alpha}^*(a, c)$  and  $\tilde{\mu}^*(a, c)$  are decreasing functions of  $a$  and  $c$ .

We have

$$COV(\tilde{\alpha}^*(\tilde{\mu}^*), \tilde{\mu}^*) = \mathbb{E}[(\tilde{\mu}^* - \mathbb{E}[\tilde{\mu}^*])(\tilde{\alpha}^*(\tilde{\mu}^*) - \mathbb{E}[\tilde{\alpha}^*(\tilde{\mu}^*)])].$$

Expanding terms, we have

$$\begin{aligned} COV(\tilde{\alpha}^*(\tilde{\mu}^*), \tilde{\mu}^*) &= \mathbb{E}[(\tilde{\mu}^* - \mathbb{E}[\tilde{\mu}^*])(\tilde{\alpha}^*(\tilde{\mu}^*) - \tilde{\alpha}^*(\mathbb{E}[\tilde{\mu}^*]))] \\ &\quad + \mathbb{E}[(\tilde{\mu}^* - \mathbb{E}[\tilde{\mu}^*])(\tilde{\alpha}^*(\mathbb{E}[\tilde{\mu}^*]) - \mathbb{E}[\tilde{\alpha}^*(\tilde{\mu}^*)])]. \end{aligned}$$

The term on the second line is zero because  $\tilde{\alpha}^*(\mathbb{E}[\tilde{\mu}^*]) - \mathbb{E}[\tilde{\alpha}^*(\tilde{\mu}^*)]$  is non-stochastic and  $\mathbb{E}[\tilde{\mu}^* - \mathbb{E}[\tilde{\mu}^*]] = 0$ . Hence, we have

$$COV(\tilde{\alpha}^*(\tilde{\mu}^*), \tilde{\mu}^*) = \mathbb{E}[(\tilde{\mu}^* - \mathbb{E}[\tilde{\mu}^*])(\tilde{\alpha}^*(\tilde{\mu}^*) - \tilde{\alpha}^*(\mathbb{E}[\tilde{\mu}^*]))] > 0.$$

The conclusion follows from the fact that  $\tilde{\alpha}^*(\tilde{\mu}^*)$  is non-decreasing, so  $\tilde{\alpha}^*(\tilde{\mu}^*) - \tilde{\alpha}^*(\mathbb{E}[\tilde{\mu}^*]) \geq 0$  if  $\tilde{\mu}^* - \mathbb{E}[\tilde{\mu}^*] \geq 0$ . Since the function  $\tilde{\alpha}^*(\tilde{\mu}^*)$  is strictly increasing on a set of positive measure, the strict inequality holds.

The proof for the covariance of  $\tilde{\mu}^*$  and  $\tilde{\sigma}^*$  is identical. In particular, because  $\mu^*(a, c)$  and  $\sigma^*(a, c)$  are comonotone, we can write  $\tilde{\sigma}^*$  as a monotonic function of  $\tilde{\mu}^*$ ,  $\tilde{\sigma}^*(\tilde{\mu}^*)$ . The remainder of the argument is then exactly as stated above. ■

## A.2 The Case of Combined Adverse Selection and Moral Hazard

**Proof of Lemma 1.1:** A simple, and incentive compatible way to exclude types  $\theta < \theta_m$  is to offer the null contract  $\beta(\theta) = \alpha(\theta) = 0$  for all types  $\theta < \theta_m$ . Assume thus that the principal offers such a scheme.

Consider now incentive compatibility. A type  $\theta$  should not have an incentive to mimic any type  $\hat{\theta}$ , so

$$\beta(\theta) + (\Delta - \Gamma)\theta\alpha(\theta)^\eta \geq \beta(\hat{\theta}) + (\Delta - \Gamma)\theta\alpha(\hat{\theta})^\eta.$$

Likewise, a type  $\hat{\theta}$  should not have an incentive to mimic any type  $\theta$ , so

$$\beta(\hat{\theta}) + (\Delta - \Gamma)\hat{\theta}\alpha(\hat{\theta})^\eta \geq \beta(\theta) + (\Delta - \Gamma)\hat{\theta}\alpha(\theta)^\eta.$$

Adding the two constraints, and rearranging, we get

$$(\Delta - \Gamma)(\theta - \hat{\theta})(\alpha(\theta)^\eta - \alpha(\hat{\theta})^\eta) \geq 0.$$

Hence  $\alpha(\theta)$  must be nondecreasing in  $\theta$ .

Let  $\omega(\theta) \equiv \max_{\hat{\theta}} \{\beta(\hat{\theta}) + (\Delta - \Gamma)\theta\alpha(\hat{\theta})^\eta\}$ . Given truthtelling at the optimum, it follows that  $\omega_\theta(\theta) = (\Delta - \Gamma)\alpha(\theta)^\eta$  almost everywhere. Imposing the participation constraint

at  $\theta_m$ , we get

$$\omega(\theta) = \int_{\underline{\theta}}^{\theta} (\Delta - \Gamma) \alpha(z)^\eta dz = \int_{\theta_m}^{\theta} (\Delta - \Gamma) \alpha(z)^\eta dz,$$

where the second equality follows from the fact that  $\alpha(\theta) = 0$  for all  $\theta < \theta_m$ . Since

$$\omega(\theta) = \beta(\theta) + (\Delta - \Gamma) \theta \alpha(\theta)^\eta,$$

we have

$$\beta(\theta) = \int_{\theta_m}^{\theta} (\Delta - \Gamma) \alpha(z)^\eta dz - (\Delta - \Gamma) \theta \alpha(\theta)^\eta \text{ for all } \theta \geq \theta_m.$$

The proof that these conditions are sufficient is standard and thus omitted.

For completeness, observe that no type  $\theta < \theta_m$  has any incentive to mimic any type  $\hat{\theta} \geq \theta_m$  by the standard reasoning. In particular, the utility such a type can get this way is

$$\begin{aligned} & \beta(\hat{\theta}) + (\Delta - \Gamma) (\theta - \hat{\theta}) \alpha(\hat{\theta})^\eta + (\Delta - \Gamma) \hat{\theta} \alpha(\hat{\theta})^\eta \\ &= \omega(\hat{\theta}) + (\Delta - \Gamma) (\theta - \hat{\theta}) \alpha(\hat{\theta})^\eta \end{aligned}$$

However, since  $\omega_\theta(\theta) = (\Delta - \Gamma) \alpha(\theta)^\eta$  and  $\omega(\theta_m) = \omega$ , we have

$$\begin{aligned} & \omega(\hat{\theta}) + (\Delta - \Gamma) (\theta - \hat{\theta}) \alpha(\hat{\theta})^\eta \\ &= \omega + \int_{\theta_m}^{\hat{\theta}} (\Delta - \Gamma) \alpha(z)^\eta dz + (\Delta - \Gamma) (\theta - \hat{\theta}) \alpha(\hat{\theta})^\eta \\ &= \omega + \int_{\theta_m}^{\hat{\theta}} (\Delta - \Gamma) (\alpha(z)^\eta - \alpha(\hat{\theta})^\eta) dz + (\Delta - \Gamma) (\theta - \theta_m) \alpha(\hat{\theta})^\eta < \omega \end{aligned}$$

where the inequality follows from the monotonicity of  $\alpha(\theta)$  and the fact that  $\theta < \theta_m$ . ■

**Proof of Lemma 1.2:** We demonstrate the following three facts:

i) for  $\theta \in [\underline{\theta}, \bar{\theta}]$ ,

$$F(\theta) = \int_x^{\min\{\bar{r}, \frac{\theta}{s}\}} G\left(\frac{\theta}{r} \middle| r\right) q(r) dr,$$

and

$$f(\theta) = \int_{\underline{r}}^{\min\{\bar{r}, \frac{\theta}{\underline{s}}\}} \frac{1}{r} g\left(\frac{\theta}{r} \mid r\right) q(r) dr,$$

where  $G(s|r)$  ( $g(s|r)$ ) is the cdf (pdf) of the conditional distribution of  $s$  given  $r$  and  $q(r)$  is the density of the marginal distribution of  $r$ ;

ii) the distribution satisfies  $f(\underline{\theta}) = 0$  and  $f(\theta) > 0$  for  $\theta > \underline{\theta}$ ;

iii) the distribution satisfies, for  $\theta > \underline{\theta}$ ,

$$\frac{\partial}{\partial \theta} \frac{1 - F(\theta)}{\theta f(\theta)} \leq 0$$

if  $g(s|r)$  satisfies

$$\frac{g_s(s|r)}{\frac{g(s|r)}{s}} \geq -1.$$

i) Consider the random variable

$$\tilde{\theta} = rs.$$

With a slight abuse of notation, let  $\theta$  denote the level that the rv  $\tilde{\theta}$  takes. Let  $h(r, s)$  denote the joint density of  $r$  and  $s$ . Hence, for  $\underline{r}\underline{s} \leq \theta \leq \bar{r}\underline{s}$ ,

$$\Pr[\tilde{\theta} \leq \theta] = \Pr[rs \leq \theta] = \int_{\underline{r}}^{\frac{\theta}{\underline{s}}} \int_{\underline{s}}^{\frac{\theta}{r}} h(r, s) ds dr, \quad (\text{A.7})$$

while for  $\bar{r}\underline{s} \leq \theta \leq \bar{r}\bar{s}$ ,

$$\Pr[\tilde{\theta} \leq \theta] = \Pr[rs \leq \theta] = \int_{\underline{r}}^{\bar{r}} \int_{\underline{s}}^{\frac{\theta}{r}} h(r, s) ds dr.$$

We treat these two cases in sequence now beginning with the former.

We can rewrite (A.7), for  $\frac{\theta}{\underline{s}} < \bar{r}$ ,

$$\int_{\underline{r}}^{\frac{\theta}{\underline{s}}} \int_{\underline{s}}^{\frac{\theta}{r}} h(r, s) ds dr = \int_{\underline{r}}^{\frac{\theta}{\underline{s}}} \int_{\underline{s}}^{\frac{\theta}{r}} g(s|r) ds q(r) dr = \int_{\underline{r}}^{\frac{\theta}{\underline{s}}} G\left(\frac{\theta}{r} \mid r\right) q(r) dr. \quad (\text{A.8})$$

The derivative of this expression with respect to  $\theta$  is

$$\frac{\partial}{\partial \theta} \left[ \int_{\underline{r}}^{\frac{\theta}{\underline{s}}} G\left(\frac{\theta}{r} \middle| r\right) q(r) dr \right] = \int_{\underline{r}}^{\frac{\theta}{\underline{s}}} \frac{1}{r} g\left(\frac{\theta}{r} \middle| r\right) q(r) dr + G\left(\frac{\theta}{\frac{\theta}{\underline{s}}} \middle| r = \frac{\theta}{\underline{s}}\right) q\left(\frac{\theta}{\underline{s}}\right).$$

Since  $G\left(\frac{\theta}{\frac{\theta}{\underline{s}}} \middle| r = \frac{\theta}{\underline{s}}\right) = 0$ , this simplifies to

$$\frac{\partial}{\partial \theta} \left[ \int_{\underline{r}}^{\frac{\theta}{\underline{s}}} G\left(\frac{\theta}{r} \middle| r\right) q(r) dr \right] = \int_{\underline{r}}^{\frac{\theta}{\underline{s}}} \frac{1}{r} g\left(\frac{\theta}{r} \middle| r\right) q(r) dr. \quad (\text{A.9})$$

For  $\frac{\theta}{\underline{s}} > \bar{r}$ , we can write

$$\Pr[\tilde{\theta} \leq \theta] = \int_{\underline{r}}^{\bar{r}} \int_{\underline{s}}^{\frac{\theta}{\bar{r}}} h(r, s) ds dr = \int_{\underline{r}}^{\bar{r}} \int_{\underline{s}}^{\frac{\theta}{\bar{r}}} g(s | r) ds q(r) dr = \int_{\underline{r}}^{\bar{r}} G\left(\frac{\theta}{r} \middle| r\right) q(r) dr. \quad (\text{A.10})$$

The derivative of this expression with respect to  $\theta$  is

$$\int_{\underline{r}}^{\bar{r}} \frac{1}{r} g\left(\frac{\theta}{r} \middle| r\right) q(r) dr \quad (\text{A.11})$$

It follows that we can write the density for all  $\theta$  as

$$f(\theta) = \int_{\underline{r}}^{\min\{\bar{r}, \frac{\theta}{\underline{s}}\}} \frac{1}{r} g\left(\frac{\theta}{r} \middle| r\right) q(r) dr,$$

and the cdf for all  $\theta$  as

$$F(\theta) = \int_{\underline{r}}^{\min\{\bar{r}, \frac{\theta}{\underline{s}}\}} G\left(\frac{\theta}{r} \middle| r\right) q(r) dr.$$

ii) Evaluating the density at

$$\underline{\theta} = \underline{r}\underline{s},$$



we obtain

$$\int_{\underline{r}}^{\bar{r}} \frac{1}{r} g\left(\frac{rs}{r} \middle| r\right) q(r) dr = 0.$$

That is, the density goes to zero at the low end. It is easy to see that the density is strictly positive for  $\theta > \underline{\theta}$ .

iii)

$$\frac{1 - F(\theta)}{\theta f(\theta)} = \frac{1 - \int_{\underline{r}}^{\min\{\bar{r}, \frac{\theta}{s}\}} G\left(\frac{\theta}{r} \middle| r\right) q(r) dr}{\int_{\underline{r}}^{\min\{\bar{r}, \frac{\theta}{s}\}} \frac{\theta}{r} g\left(\frac{\theta}{r} \middle| r\right) q(r) dr}.$$

Differentiating for  $\theta > \bar{r}s$ , we find that  $\frac{\partial}{\partial \theta} \frac{1-F(\theta)}{\theta f(\theta)}$  is proportional to

$$-\theta \left( \int_{\underline{r}}^{\bar{r}} \frac{1}{r} g\left(\frac{\theta}{r} \middle| r\right) q(r) dr \right)^2 - \left( \int_{\underline{r}}^{\bar{r}} \left( g\left(\frac{\theta}{r} \middle| r\right) + \frac{\theta}{r} g_s\left(\frac{\theta}{r} \middle| r\right) \right) \frac{q(r)}{r} dr \right) \cdot \left( 1 - \int_{\underline{r}}^{\bar{r}} G\left(\frac{\theta}{r} \middle| r\right) q(r) dr \right).$$

The expression is negative if  $g(s|r) + sg_s(s|r) \geq 0$ , which is satisfied if the elasticity of the density is larger than minus unity,

$$\frac{g_s(s|r)}{\frac{g(s|r)}{s}} \geq -1.$$

Differentiating for  $\theta < \bar{r}\underline{s}$ , we find that  $\frac{\partial}{\partial \theta} \frac{1-F(\theta)}{\theta f(\theta)}$  is proportional to

$$\begin{aligned} & - \left( \int_{\frac{\underline{s}}{r}}^{\frac{\theta}{\underline{s}}} \frac{1}{r} g \left( \frac{\theta}{r} \middle| r \right) q(r) dr + G \left( \frac{\underline{s}}{r} \middle| \frac{\theta}{\underline{s}} \right) q \left( \frac{\theta}{\underline{s}} \right) \frac{1}{\underline{s}} \right) \int_{\frac{\underline{s}}{r}}^{\frac{\theta}{\underline{s}}} \frac{\theta}{r} g \left( \frac{\theta}{r} \middle| r \right) q(r) dr \\ & - \left( \int_{\frac{\underline{s}}{r}}^{\frac{\theta}{\underline{s}}} \frac{1}{r} g \left( \frac{\theta}{r} \middle| r \right) q(r) dr + \int_{\frac{\underline{s}}{r}}^{\frac{\theta}{\underline{s}}} \frac{\theta}{r^2} g_s \left( \frac{\theta}{r} \middle| r \right) q(r) dr + \underline{s} g \left( \frac{\underline{s}}{r} \middle| \frac{\theta}{\underline{s}} \right) q \left( \frac{\theta}{\underline{s}} \right) \frac{1}{\underline{s}} \right) \\ & \left( 1 - \int_{\frac{\underline{s}}{r}}^{\frac{\theta}{\underline{s}}} G \left( \frac{\theta}{r} \middle| r \right) q(r) dr \right). \end{aligned}$$

Since  $G \left( \frac{\underline{s}}{r} \middle| \frac{\theta}{\underline{s}} \right) = 0$ , the same condition on the conditional density of  $r$  implies the result. ■

**Proof of Proposition 1.4:** Using Leibniz' rule for differentiation of integrals, we find

$$\begin{aligned} V'(\theta_m) = & - \left\{ \left( \theta_m \Delta \alpha(\theta_m)^{\eta-1} - \Gamma \theta_m \alpha(\theta_m)^\eta \right) f(\theta_m) - (\Delta - \Gamma) \alpha(\theta_m)^\eta (1 - F(\theta_m)) \right\} \\ & + \omega f(\theta_m). \end{aligned}$$

Recall the first-order condition for the optimal  $\alpha$  :

$$\left( (\eta - 1) \theta \Delta \alpha(\theta)^{\eta-2} - \Gamma \theta \eta \alpha(\theta)^{\eta-1} \right) f(\theta) - (\Delta - \Gamma) \eta \alpha(\theta)^{\eta-1} (1 - F(\theta)) = 0.$$

Multiplying by  $\alpha(\theta)$  and rearranging, we find that

$$\left( \theta \Delta \alpha(\theta)^{\eta-1} - \Gamma \theta \alpha(\theta)^\eta \right) f(\theta) - (\Delta - \Gamma) \alpha(\theta)^\eta (1 - F(\theta)) = \theta \frac{\Delta}{\eta} \alpha(\theta)^{\eta-1} f(\theta).$$

Noting that  $\mu(\theta) = \theta \Delta \alpha(\theta)^{\eta-1}$ , the result follows. ■

**Proof of proposition 1.5:** Part i): Note that  $\alpha^*(\theta)$  and  $\mu^*(\theta) = \Delta \theta \alpha(\theta)^{\eta-1}$  are nondecreasing functions of  $\theta$ . Applying the same argument as in the proof of proposition 1.2 the result follows immediately.

Part ii): Since  $\alpha^*(\theta)$  is nondecreasing in  $\theta$ , the proof of proposition 1.2 can be extended to the present case if  $\sigma^*$  is nondecreasing in its arguments  $r$  and  $s$  and if these random variables are associated.<sup>1</sup>

<sup>1</sup>In principle, one could apply the same logic to the case where  $\sigma$  is decreasing in its arguments to conclude

We now give conditions for the monotonicity properties of  $\sigma(\alpha, \mathbf{t})^*$ . Recall that  $r \equiv a^{\frac{-\delta}{2(1-\lambda)-\delta}}$  and  $s \equiv c^{\frac{-2\lambda}{2(1-\lambda)-\delta}}$  and that

$$\begin{aligned}\sigma(\alpha, \mathbf{t})^* &= \Lambda a^{\frac{\lambda-1}{2(1-\lambda)-\delta}} c^{\frac{-\lambda}{2(1-\lambda)-\delta}} \alpha(\theta)^{\frac{2\lambda-1}{2(1-\lambda)-\delta}} \\ &= \Lambda r^{\frac{1-\lambda}{\delta}} s^{\frac{1}{2}} \alpha(rs)^{\frac{2\lambda-1}{2(1-\lambda)-\delta}}.\end{aligned}$$

We have

$$\frac{\partial \sigma(\alpha, \mathbf{t})^*}{\partial r} = \frac{1-\lambda}{\delta} \frac{\sigma(\alpha, \mathbf{t})^*}{r} + \frac{2\lambda-1}{2(1-\lambda)-\delta} \frac{\sigma(\alpha, \mathbf{t})^*}{\alpha} \frac{\partial \alpha(\theta)}{\partial \theta} s,$$

so  $\frac{\partial \sigma(\alpha, \mathbf{t})^*}{\partial r} \geq 0 (\leq 0)$  if and only if

$$\alpha(\theta) \geq (\leq) \frac{\delta(1-2\lambda)}{(1-\lambda)(2(1-\lambda)-\delta)} \frac{\partial \alpha(\theta)}{\partial \theta} \theta.$$

Similarly, we have

$$\frac{\partial \sigma(\alpha, \mathbf{t})^*}{\partial s} = \frac{1}{2} \frac{\sigma(\alpha, \mathbf{t})^*}{s} + \frac{2\lambda-1}{2(1-\lambda)-\delta} \frac{\sigma(\alpha, \mathbf{t})^*}{\alpha(\theta)} \frac{\partial \alpha(\theta)}{\partial \theta} r,$$

so  $\frac{\partial \sigma(\alpha, \mathbf{t})^*}{\partial s} \geq 0 (\leq 0)$  if and only if

$$\alpha(\theta) \geq (\leq) \frac{2(1-2\lambda)}{2(1-\lambda)-\delta} \frac{\partial \alpha(\theta)}{\partial \theta} \theta.$$

Noting that  $2 \geq \frac{\delta}{1-\lambda}$  by assumption,  $\sigma(\alpha, \mathbf{t})$  is increasing in both arguments iff

$$\alpha(\theta) \geq \frac{2(1-2\lambda)}{2(1-\lambda)-\delta} \frac{\partial \alpha(\theta)}{\partial \theta} \theta.$$

Likewise,  $\sigma(\alpha, \mathbf{t})$  is decreasing in both arguments iff

$$\alpha(\theta) \leq \frac{\delta(1-2\lambda)}{(1-\lambda)(2(1-\lambda)-\delta)} \frac{\partial \alpha(\theta)}{\partial \theta} \theta.$$

---

that the covariance between  $\tilde{\alpha}^*$  and  $-\tilde{\sigma}^*$  becomes positive (and hence the covariance between  $\tilde{\alpha}^*$  and  $\tilde{\sigma}^*$  negative); however, we have not been able to find a meaningful sufficient condition on the distribution that ensure this.

Differentiating the optimal  $\alpha$  with respect to  $\theta$ , we obtain

$$\begin{aligned}
\alpha_\theta(\theta) &= -\frac{\eta-1}{\eta} \Delta \left( \Gamma + (\Delta - \Gamma) \frac{1-F(\theta)}{\theta f(\theta)} \right)^{-2} (\Delta - \Gamma) \frac{\partial}{\partial \theta} \frac{1-F(\theta)}{\theta f(\theta)} \\
&= -\alpha(\theta) \frac{\frac{\partial}{\partial \theta} \frac{1-F(\theta)}{\theta f(\theta)}}{\frac{\Gamma}{\Delta - \Gamma} + \frac{1-F(\theta)}{\theta f(\theta)}} \\
&= -\alpha(\theta) \frac{\frac{\partial}{\partial \theta} \frac{1-F(\theta)}{\theta f(\theta)}}{\frac{\delta+2\lambda}{2(1-\lambda)-\delta} + \frac{1-F(\theta)}{\theta f(\theta)}}
\end{aligned}$$

Therefore,  $\sigma$  is increasing in both arguments iff

$$\alpha(\theta) \geq -\alpha(\theta) \frac{2(1-2\lambda)}{2(1-\lambda)-\delta} \frac{\frac{\partial}{\partial \theta} \frac{1-F(\theta)}{\theta f(\theta)}}{\frac{\delta+2\lambda}{2(1-\lambda)-\delta} + \frac{1-F(\theta)}{\theta f(\theta)}} \theta$$

Since  $\delta \leq 2\lambda$  implies that  $\frac{2(1-2\lambda)}{2(1-\lambda)-\delta} \leq 1$ , a sufficient condition is

$$-\frac{\partial}{\partial \theta} \frac{1-F(\theta)}{\theta f(\theta)} \theta \leq \frac{\delta+2\lambda}{2(1-\lambda)-\delta} + \frac{1-F(\theta)}{\theta f(\theta)}$$

which is equivalent to

$$-\frac{\partial}{\partial \theta} \frac{1-F(\theta)}{f(\theta)} \leq \frac{\delta+2\lambda}{2(1-\lambda)-\delta},$$

the condition given in the proposition. ■





# Appendix B

## Subgroup Deliberation and Voting

### B.1 The Model

#### Proof of Lemma 2.3:

**Step 1** In a reactive SSDE, two types of individual deviations must be prevented. The first type involves a deviation at the voting stage following a truthful announcement at the communication stage. The second type of deviation involves lying at the communication stage.

**Step 2** We here prove Point a), corresponding to the set of type 2 reactive SSDEs. We first show that the condition given in Point a) is *sufficient* to ensure that none of the above mentioned two types of deviations is strictly advantageous to a juror of type  $j$ . Assume thus that the condition of Point a) is satisfied. Regarding the first type of mentioned deviation, the threshold adopted by each Subgroup is ex post optimal at the voting stage, conditional on the locally pooled information and assuming individual pivotality, i.e. assuming that that the other Subgroup votes for conviction. We now examine the second type of deviation. Note that misreporting a  $g$ -signal as an  $i$ -signal is either inconsequential or adversely triggers an acquittal given a Subgroup signal profile where the deviating juror would have favoured a conviction. This can thus not be strictly advantageous to a juror. Instead, misreporting an  $i$ -signal as a  $g$ -signal is always without consequence on the final decision, as a juror can always block a conviction triggered by his lie if he realizes that he favours acquittal, given remaining Subgroup members' signals.

We now show that the condition stated in Point a) is *necessary* to ensure that none of the two types of deviations mentioned in step 1 is strictly advantageous to a juror of type  $j$ .

Suppose that thus that the condition is not satisfied. Suppose that  $t_j$  is larger than specified by the condition, given  $t_{-j}$ . Then a juror of preference type  $j$  has a strict incentive to announce an  $i$ -signal as a  $g$ -signal and subsequently vote on the basis of the known signal profile of his Subgroup and the assumption that the other Subgroup convicts. Suppose now instead that  $t_j$  is smaller than specified by the condition, given  $t_{-j}$ . Then a juror of preference type  $j$  has a strict incentive to announce a  $g$ -signal as an  $i$ -signal and subsequently vote on the basis of the known signal profile of his Subgroup and the assumption that the other Subgroup convicts.

**Step 3** We now prove Point b), corresponding to the set of type 1 reactive SSDEs. The analysis of condition (2.9) for type  $j$  follows the exact same steps as in Point a). We now examine condition (2.10), which applies to the type that always convicts independently of the its Subgroup signal profile. Note first that a juror of type  $-j$  must be willing to convict no matter what signal profile is revealed at the communication stage, which requires (2.10) to hold. This proves that (2.10) is *necessary*. We now show that condition (2.10) is *sufficient* to ensure no strict incentive to deviate for type  $-j$ . An individual of type  $-j$  recognizes that his announced signal is inconsequential for the voting behavior of his Subgroup and thus has no incentive to deviate from truthtelling. As to the voting stage, conviction is always ex post optimal, assuming individual pivotality, i.e. assuming that that the other Subgroup votes for conviction. It follows that a type  $-j$  has no strict incentive to deviate at the voting stage.

**Step 4** In the next steps, we show that our characterization of the set of reactive SSDEs generalizes to a larger set of voting rules. Let  $R$  be the minimal number of conviction votes required for a conviction decision and assume that  $R > \{n_H, n_D\}$ . Two key aspects deserve mention. First, assuming  $R > \{n_H, n_D\}$  means that individual pivotality, either in communicating or in voting, implies that the Subgroup to which one does not belong votes for conviction. This replicates the case of Unanimity. A second key aspect is that abandoning Unanimity implies that an individual can now not single handedly veto a conviction anymore. Accordingly, deviating to announcing a  $g$ -signal when holding an  $i$ -signal is now risky, in the sense that one cannot simply veto an undesirable collective conviction vote triggered by such a deviation. We now show that the necessary and sufficient conditions given for the case of Unanimity, whether in Point a) or Point b), extend to this more general case.

**Step 5** We first look at the set of type 2 reactive SSDEs. We first show that the condition of Point a) is *sufficient* to ensure that none of the two types of deviations identified in step 1 is strictly advantageous. Assume thus that condition of Point a) is respected. Regarding



the first type of mentioned deviation, the threshold adopted by each Subgroup is ex post optimal at the voting stage, conditional on the locally pooled information and assuming individual pivotality, i.e. assuming that the other Subgroup votes for conviction. We now examine the second type of deviation. Note that misreporting a  $g$ -signal as an  $i$ -signal is either inconsequential or adversely triggers an acquittal given a signal profile where the deviating juror would have favoured a conviction. This can thus not be strictly advantageous to a juror. Instead, misreporting an  $i$ -signal as a  $g$ -signal is either inconsequential or adversely triggers a conviction given a signal profile where the deviating juror would have favoured an acquittal. This can thus not be strictly advantageous to a juror.

We now show that the condition given in Point a) is *necessary* to ensure that none of the two types of deviations mentioned in step 1 is strictly advantageous. Suppose thus that the condition is not satisfied. Suppose that  $t_j$  is larger than specified by the condition, given  $t_{-j}$ . Then a juror of preference type  $j$  has a strict incentive to announce an  $i$ -signal as a  $g$ -signal and subsequently vote on the basis of the known signal profile of his Subgroup and the assumption that the other Subgroup convicts. Suppose that instead  $t_j$  is smaller than specified by the condition, given  $t_{-j}$ . Then a juror of preference type  $j$  has a strict incentive to announce a  $g$ -signal as an  $i$ -signal and subsequently vote on the basis of the known signal profile of his Subgroup and the assumption that the other Subgroup convicts.

**Step 6** We now examine the set of type 1 reactive SSDEs. The analysis of (2.9) for type  $j$  follows the exact same steps as the analysis of type 2 reactive SSDEs. The analysis of (2.10), corresponding to type  $-j$ , is identical to that given in step 3 and thus not repeated. ■

### A further Lemma on reactive SSDEs

The following lemma states in close form the existence conditions for a type 2 reactive SSDE.

#### Lemma B.1 *SSDEs*.

$(t_H, t_D)$  constitutes a type 2 reactive SSDE iff,  $\forall j \in \{H, D\}$ , it holds that  $t_j \in \{1, \dots, n_j\}$  and

$$\frac{F(p, q_j) + n_j + \Omega(p, t_{-j}, n_{-j})}{2} < t_j \leq \frac{F(p, q_j) + n_j + \Omega(p, t_{-j}, n_{-j}) + 2}{2}, \quad (\text{B.1})$$

where

$$F(p, q) \equiv \frac{\ln\left(\frac{q}{1-q}\right)}{\ln\left(\frac{p}{1-p}\right)} \text{ and } \Omega(p, k, n) \equiv \frac{\ln\left(\frac{\sum_{x \geq k}^n B(1-p, x, n)}{\sum_{x \geq k}^n B(p, x, n)}\right)}{\ln\left(\frac{p}{1-p}\right)}. \quad (\text{B.2})$$

**Proof of Lemma B.1:** Note that  $(t_H, t_D)$  constitutes a type 2 reactive SSDE iff,  $\forall j \in \{H, D\}$ , it holds that  $t_j \in \{1, \dots, n_j\}$  and the following two inequalities simultaneously hold:

$$\frac{B(p, t_j - 1, n_j) \left( \sum_{x \geq t_j}^{n-j} B(p, x, n-j) \right)}{B(p, t_j - 1, n_j) \left( \sum_{x \geq t_j}^{n-j} B(p, x, n-j) \right) + B(1-p, t_j - 1, n_j) \left( \sum_{x \geq t_j}^{n-j} B(1-p, x, n-j) \right)} < q_j \quad (\text{B.3})$$

and

$$q_j \leq \frac{B(p, t_j, n_j) \left( \sum_{x \geq t_j}^{n-j} B(p, x, n-j) \right)}{B(p, t_j, n_j) \left( \sum_{x \geq t_j}^{n-j} B(p, x, n-j) \right) + B(1-p, t_j, n_j) \left( \sum_{x \geq t_j}^{n-j} B(1-p, x, n-j) \right)}. \quad (\text{B.4})$$

Now, note that (B.3) can be rewritten as follows:

$$\begin{aligned} & (1 - q_j) p^{t_j - 1} (1 - p)^{n_j - t_j + 1} \left( \sum_{x \geq t_j}^{n-j} B(p, x, n-j) \right) \\ & < q_j (1 - p)^{t_j - 1} p^{n_j - t_j + 1} \left( \sum_{x \geq t_j}^{n-j} B(1-p, x, n-j) \right). \end{aligned} \quad (\text{B.5})$$

Applying the ln-transformation to both sides of (B.5), the above inequality can then be rewritten as follows:

$$\frac{\ln\left(\frac{q_j}{1-q_j}\right)}{2 \ln\left(\frac{p}{1-p}\right)} + \frac{\ln\left(\frac{\sum_{x \geq t_j}^{n-j} B(1-p, x, n-j)}{\sum_{x \geq t_j}^{n-j} B(p, x, n-j)}\right)}{2 \ln\left(\frac{p}{1-p}\right)} + \frac{n_j}{2} < t_j. \quad (\text{B.6})$$

One can perform a similar transformation for (B.4). One obtains an inequality stating that  $t_j$  is weakly smaller than the LHS expression in (B.6) plus one. ■

## B.2 Positive Analysis

**Proof of Lemma 2.4 (reactive SNDEs):**

**Step 1** We first analyze the set of reactive SNDEs in which both preference types condition their play on their information. Note that a given preference type cannot mix after both  $i$ - and  $g$ -signals (see condition (2.4)). Within this subclass of equilibria, there are altogether

nine possible symmetric voting profiles which are listed and numbered in the Table below. Letters  $x, y \in (0, 1)$  are used to denote mixing probabilities.

	$\sigma_g^H, \sigma_i^H$	$\sigma_g^D, \sigma_i^D$		$\sigma_g^H, \sigma_i^H$	$\sigma_g^D, \sigma_i^D$		$\sigma_g^H, \sigma_i^H$	$\sigma_g^D, \sigma_i^D$
<b>1</b>	1, 0	1, 0	<b>4</b>	$x, 0$	1, 0	<b>7</b>	$x, 0$	1, $y$
<b>2</b>	1, 0	$x, 0$	<b>5</b>	1, $x$	1, 0	<b>8</b>	1, $x$	$y, 0$
<b>3</b>	1, 0	1, $x$	<b>6</b>	$x, 0$	$y, 0$	<b>9</b>	1, $x$	1, $y$

We show that none of the above nine strategy profiles constitutes an equilibrium. Equilibrium 1 trivially never exists when  $m > 1$ . Equilibria 2,4 and 6 do not exist under the assumption that  $q_D < \beta(p, n, n)$  given that they require either  $q_D = \beta(p, n, n)$  or  $q_H = \beta(p, n, n)$  (recall  $q_H < q_D$ ). Recall in what follows that  $piv_j$  stands for the event in which a juror of preference type  $j$  is pivotal, i.e. all remaining jurors vote for conviction. Equilibria 3,7 and 9 imply (B.7) and (B.8), as given below.

$$\begin{aligned}
q_D &= P(G \mid i, piv_D) && \text{(B.7)} \\
&= \frac{(1-p) [p\sigma_g^D + (1-p)\sigma_i^D]^{n_D-1} [p\sigma_g^H + (1-p)\sigma_i^H]^{n_H}}{\left( (1-p) [p\sigma_g^D + (1-p)\sigma_i^D]^{n_D-1} [p\sigma_g^H + (1-p)\sigma_i^H]^{n_H} \right.} \\
&\quad \left. + p [(1-p)\sigma_g^D + p\sigma_i^D]^{n_D-1} [(1-p)\sigma_g^H + p\sigma_i^H]^{n_H} \right) \\
&\leq \frac{(1-p)p [p\sigma_g^D + (1-p)\sigma_i^D]^{n_D-1} [p\sigma_g^H + (1-p)\sigma_i^H]^{n_H-1}}{\left( (1-p)p [p\sigma_g^D + (1-p)\sigma_i^D]^{n_D-1} [p\sigma_g^H + (1-p)\sigma_i^H]^{n_H-1} \right.} \\
&\quad \left. + p(1-p) [(1-p)\sigma_g^D + p\sigma_i^D]^{n_D-1} [(1-p)\sigma_g^H + p\sigma_i^H]^{n_H-1} \right) \\
&= \frac{pF_p^1}{pF_p^1 + (1-p)F_{1-p}^1} \equiv \bar{P}_1,
\end{aligned}$$

$$\begin{aligned}
q_H &\geq P(G \mid i, piv_H) & (B.8) \\
&= \frac{(1-p) [p\sigma_g^D + (1-p)\sigma_i^D]^{n_D} [p\sigma_g^H + (1-p)\sigma_i^H]^{n_H-1}}{\left( (1-p) [p\sigma_g^D + (1-p)\sigma_i^D]^{n_D} [p\sigma_g^H + (1-p)\sigma_i^H]^{n_H-1} \right. \\
&\quad \left. + p [(1-p)\sigma_g^D + p\sigma_i^D]^{n_D} [(1-p)\sigma_g^H + p\sigma_i^H]^{n_H-1} \right)} \\
&> \frac{(1-p)^2 [p\sigma_g^D + (1-p)\sigma_i^D]^{n_D-1} [p\sigma_g^H + (1-p)\sigma_i^H]^{n_H-1}}{\left( (1-p)^2 [p\sigma_g^D + (1-p)\sigma_i^D]^{n_D-1} [p\sigma_g^H + (1-p)\sigma_i^H]^{n_H-1} \right. \\
&\quad \left. + p^2 [(1-p)\sigma_g^D + p\sigma_i^D]^{n_D-1} [(1-p)\sigma_g^H + p\sigma_i^H]^{n_H-1} \right)} \\
&= \frac{(1-p)F_p^1}{(1-p)F_p^1 + pF_{1-p}^1} \equiv \underline{P_1},
\end{aligned}$$

where

$$F_r^1 \equiv (1-r) [r\sigma_g^D + (1-r)\sigma_i^D]^{n_D-1} [r\sigma_g^H + (1-r)\sigma_i^H]^{n_H-1}, \quad r \in \{p, (1-p)\}.$$

Now, using the fact that for any positive constants  $A, B, C, D$ ,  $\frac{A}{A+B} \leq \frac{C}{C+D} \Leftrightarrow \frac{A}{B} \leq \frac{C}{D}$ , note that there exists a positive integer  $T$  s.t.

$$\frac{B(p, T-1, n)}{B(1-p, T-1, n)} = \frac{p^{T-1}(1-p)^{n-T+1}}{(1-p)^{T-1}p^{n-T+1}} \leq \frac{(1-p)F_p^1}{pF_{1-p}^1} \leq \frac{p^T(1-p)^{n-T}}{(1-p)^T p^{n-T}} = \frac{B(p, T, n)}{B(1-p, T, n)} \quad (B.9)$$

and (multiplying all expressions by  $\frac{p^2}{(1-p)^2}$ )

$$\frac{B(p, T, n)}{B(1-p, T, n)} = \frac{p^T(1-p)^{n-T}}{(1-p)^T p^{n-T}} \leq \frac{pF_p^1}{(1-p)F_{1-p}^1} \leq \frac{p^{T+1}(1-p)^{n-T-1}}{(1-p)^{T+1}p^{n-T-1}} = \frac{B(p, T+1, n)}{B(1-p, T+1, n)}. \quad (B.10)$$

Summarizing, inequalities (B.7) and (B.8) thus imply that there exists a positive integer  $T$  s.t.:

$$\beta(p, T-1, n) \leq \underline{P_1} \leq q_H < q_D \leq \overline{P_1} \leq \beta(p, T+1, n). \quad (B.11)$$

The inequality relation (B.11) however means that  $m \leq 1$  if equilibrium 3,7 or 9 exist. But we have assumed  $m > 1$ . As to equilibria 5 and 8, note that they imply that the following

two conditions (B.12) and (B.13) hold:

$$\begin{aligned}
q_H &= P(G \mid i, piv_H) & (B.12) \\
&= \frac{(1-p)[p]^{n_D} [p\sigma_g^H + (1-p)\sigma_i^H]^{n_H-1}}{\left( (1-p)[p]^{n_D} [p\sigma_g^H + (1-p)\sigma_i^H]^{n_H-1} \right. \\
&\quad \left. + p[1-p]^{n_D} [(1-p)\sigma_g^H + p\sigma_i^H]^{n_H-1} \right)} \\
&= \frac{(1-p)F_p^2}{(1-p)F_p^2 + pF_{1-p}^2} \equiv \underline{P}_2,
\end{aligned}$$

$$\begin{aligned}
q_D &\leq P(G \mid g, piv_D) & (B.13) \\
&= \frac{[p]^{n_D} [p\sigma_g^H + (1-p)\sigma_i^H]^{n_H}}{\left( [p]^{n_D} [p\sigma_g^H + (1-p)\sigma_i^H]^{n_H} \right. \\
&\quad \left. + [(1-p)]^{n_D} [(1-p)\sigma_g^H + p\sigma_i^H]^{n_H} \right)} \\
&< \frac{p[p]^{n_D} [p\sigma_g^H + (1-p)\sigma_i^H]^{n_H-1}}{\left( p[p]^{n_D} [p\sigma_g^H + (1-p)\sigma_i^H]^{n_H-1} \right. \\
&\quad \left. + (1-p)[1-p]^{n_D} [(1-p)\sigma_g^H + p\sigma_i^H]^{n_H-1} \right)} \\
&= \frac{pF_p^2}{pF_p^2 + (1-p)F_{1-p}^2} \equiv \overline{P}_2,
\end{aligned}$$

where

$$F_r^2 \equiv [r]^{n_D} [r\sigma_g^H + (1-r)\sigma_i^H]^{n_H-1}, \quad r \in \{p, (1-p)\}.$$

The inequalities (B.12) and (B.13) imply that there exists a positive integer  $T$  s.t.:

$$\beta(p, T-1, n) \leq \underline{P}_2 = q_H < q_D < \overline{P}_2 \leq \beta(p, T+1, n). \quad (B.14)$$

Now, note that (B.14) means that  $m \leq 1$  if equilibrium 5 or 8 exists. But we have assumed  $m > 1$ . To summarize Step 1, we have now shown that none of the nine possible reactive SND voting profiles in which both types condition their play on their information (as listed in the Table above) ever constitutes an equilibrium.

**Step 2** The next steps examine the set of putative reactive SNDEs in which at least one of the two preference types plays ( $\sigma_g = 1, \sigma_i = 1$ ) while the other type conditions its play on its information. Here, altogether six profiles need to be considered, depending on the nature

of the strategy,  $(\sigma_g = 1, \sigma_i = 0)$  or  $(\sigma_g = 1, \sigma_i = x)$  or  $(\sigma_g = y, \sigma_i = 0)$ ,  $0 < x, y < 1$ , played by the preference type that conditions its play on its signal as well as on the identity of the concerned preference type. Step 3 deals with the set of putative equilibria in which the hawks condition their play on their information while doves play  $(\sigma_g^D = 1, \sigma_i^D = 1)$ . We show that this set is empty. Step 4 examines equilibria in which the doves condition play on their signals while the hawks play  $(\sigma_g^H = 1, \sigma_i^H = 1)$ .

**Step 3** We here examine strategy profiles in which the hawks condition their play on their signal while the doves play  $(\sigma_g^D = 1, \sigma_i^D = 1)$ . In such an equilibrium it must be the case that:

$$P(G \mid i, piv_H) \leq q_H \leq P(G \mid g, piv_H), \quad (\text{B.15})$$

$$q_D \leq P(G \mid i, piv_D) < P(G \mid g, piv_D). \quad (\text{B.16})$$

Now, note however that:

$$\begin{aligned} P(G \mid i, piv_H) &= \frac{(1-p) [p\sigma_g^H + (1-p)\sigma_i^H]^{n_H-1}}{(1-p) [p\sigma_g^H + (1-p)\sigma_i^H]^{n_H-1} + p [(1-p)\sigma_g^H + p\sigma_i^H]^{n_H-1}} \quad (\text{B.17}) \\ &\geq \frac{(1-p)^2 [p\sigma_g^H + (1-p)\sigma_i^H]^{n_H-1}}{(1-p)^2 [p\sigma_g^H + (1-p)\sigma_i^H]^{n_H-1} + p^2 [(1-p)\sigma_g^H + p\sigma_i^H]^{n_H-1}} \\ &= \frac{(1-p)F_p^3}{(1-p)F_p^3 + pF_{1-p}^3} \equiv \underline{P}_3, \end{aligned}$$

$$\begin{aligned} P(G \mid i, piv_D) &= \frac{(1-p) [p\sigma_g^H + (1-p)\sigma_i^H]^{n_H}}{(1-p) [p\sigma_g^H + (1-p)\sigma_i^H]^{n_H} + p [(1-p)\sigma_g^H + p\sigma_i^H]^{n_H}} \quad (\text{B.18}) \\ &\leq \frac{(1-p)p [p\sigma_g^H + (1-p)\sigma_i^H]^{n_H-1}}{(1-p)p [p\sigma_g^H + (1-p)\sigma_i^H]^{n_H-1} + p(1-p) [(1-p)\sigma_g^H + p\sigma_i^H]^{n_H-1}} \\ &= \frac{pF_p^3}{pF_p^3 + (1-p)F_{1-p}^3} \equiv \overline{P}_3, \end{aligned}$$

, where

$$F_r^3 := (1-r) [r\sigma_g^H + (1-r)\sigma_i^H]^{n_H-1}, \quad r \in \{p, (1-p)\}.$$

Now, (B.17) and (B.18) imply that there exists a positive integer  $T$  s.t.:

$$\beta(p, T-1, n) \leq \underline{P}_3 \leq q_H < q_D \leq \overline{P}_3 \leq \beta(p, T+1, n). \quad (\text{B.19})$$

This in turn means that  $m \leq 1$ . We have however assumed  $m > 1$ . Therefore this type of equilibria does not exist.

**Step 4** We now examine equilibria in which the doves condition play on their signals while the hawks play ( $\sigma_g^H = 1, \sigma_i^H = 1$ ). There are a priori three such candidates. The first candidate is the equilibrium given by ( $\sigma_g^H = 1, \sigma_i^H = 1, \sigma_g^D = x, \sigma_i^D = 0$ ), for  $0 < x < 1$ . However, it exists iff  $q_D = \beta(p, n_D, n_D)$ , which is never true by assumption. The second candidate is the putative equilibrium A given by ( $\sigma_g^H = 1, \sigma_i^H = 1, \sigma_g^D = 1, \sigma_i^D = 0$ ). The third candidate is the putative equilibrium B given by ( $\sigma_g^H = 1, \sigma_i^H = 1, \sigma_g^D = 1, \sigma_i^D = y$ ), for  $0 < y < 1$ . We show that either equilibrium A or B (never both) exists for any  $q_D \in ((1-p), \beta(p, n_D, n_D))$ . Equilibrium A trivially exists iff  $\beta(p, n_D - 1, n_D) < q_D < \beta(p, n_D, n_D)$ . As to equilibrium B, note that  $y$  satisfies:

$$q_D = \frac{(1-p)[p + (1-p)y]^{n_D-1}}{(1-p)[p + (1-p)y]^{n_D-1} + p[1-p+py]^{n_D-1}}, \quad (\text{B.20})$$

so that, recalling explicitly the dependence of  $y$  on  $p, q_D$  and  $n_D$ ,

$$y(p, q_D, n_D) = \frac{\left(\frac{(1-q_D)(1-p)}{q_D p}\right)^{\frac{1}{n_D-1}} p - (1-p)}{p - \left(\frac{(1-q_D)(1-p)}{q_D p}\right)^{\frac{1}{n_D-1}} (1-p)}. \quad (\text{B.21})$$

Now, note that  $y(p, 1-p, n_D) = 1$ ,  $y(p, \beta(p, n_D - 1, n_D), n_D) = 0$  and

$$\begin{aligned} & \frac{\partial y(p, q_D, n_D)}{\partial q_D} \\ = & \frac{1}{pq_D^2 (n_D - 1) \left( p - \left(\frac{1}{pq_D} (p-1)(q_D-1)\right)^{\frac{1}{n_D-1}} + p \left(\frac{1}{pq_D} (p-1)(q_D-1)\right)^{\frac{1}{n_D-1}} \right)^2} \\ & \times \frac{2p^2 - 3p + 1}{\left(\frac{1}{pq_D} (p-1)(q_D-1)\right)^{\frac{1}{n_D-1} (n_D-2)}} \\ < & 0. \end{aligned}$$

It follows that equilibrium B exists iff  $1-p < q_D < \beta(p, n_D - 1, n_D)$ . ■

### Proof of Lemma 2.5 (reactive SPDEs):

**Step 1** Suppose a reactive SPDE in which hawks truthfully reveal their signals and doves babble. We know from Lemma 2.3 that such an equilibrium exists iff there is a  $t_H \in \{1, \dots, n_H\}$

s.t.  $\beta(p, t_H - 1, n_H) < q_H \leq \beta(p, t_H, n_H)$  and  $q_D \leq \beta(p, t_H, n_H + 1)$ . However, given our assumption that  $m > 1$ , there by definition exists no such  $t_H$ .

**Step 2** Suppose now a reactive SPDE in which doves truthfully reveal their signals and hawks babble. Given our assumption on  $q_D$ , there exists a (unique)  $t_D^* \in \{1, \dots, n_D\}$  s.t.  $\beta(p, t_D^* - 1, n_D) < q_D \leq \beta(p, t_D^*, n_D)$ . Furthermore, we know that  $q_H \leq \beta(p, t_D^*, n_D + 1)$  given our assumption that  $m > 1$ . It follows from Lemma 2.3 that there exists a unique SPDE in which doves truthfully communicate while hawks babble. ■

**Proof of Lemma 2.6 (reactive SSDEs):**

**Point a)** Note first that there exists a type 2 reactive SSDE if :

$$P\left(G \mid |g|_{-j} \geq T_j^{n_j}, |g|_j = 0\right) < q_j \leq \beta(p, n_j, n_j), \quad \forall j \in \{H, D\}. \quad (\text{B.22})$$

Note that there exists a type 1 reactive SSDE given by  $t_j \in \{1, \dots, n_j\}$  and  $t_{-j} = 0$  iff:

$$\left(\beta(p, 0, n_j) < q_j \leq \beta(p, n_j, n_j)\right) \cap \left(q_j \leq P\left(G \mid |g|_j \geq T_j^{n_j}, |g|_{-j} = 0\right)\right). \quad (\text{B.23})$$

Clearly, using together conditions (B.22) and (B.23), there always exists some reactive SSDE given our assumptions on  $q_H$  and  $q_D$ . Indeed, if  $\beta(p, 0, n_H) < q_H < \beta(p, n_H, n_H)$  and  $\beta(p, 0, n_D) < q_D < \beta(p, n_D, n_D)$ , then either (B.22) is true or (B.23) is true for some  $j \in \{H, D\}$ . Note finally that conditions (B.22) and (B.23) do not prohibit the simultaneous existence of a type 1 reactive SSDE and a type 2 reactive SSDE.

Note that there may exist multiple reactive SSDEs. We prove this by an example. Suppose  $n_H = 6, n_D = 8, q_H = 0.7, q_D = 0.9$  and  $p = 0.83$ . For these parameters, it is readily checked that there exist two type 2 reactive SSDEs given by respectively  $(t_H = 3, t_D = 4)$  and  $(t_H = 2, t_D = 5)$ .

**Point b)** Using the conditions given in Lemma B.1 in Appendix B.1, call  $t_i^{BR}(t_j)$  the unique best response threshold of Subgroup  $i$  to the threshold  $t_j$  of Subgroup  $j$ , as defined in (B.1). Note that either  $t_i^{BR}(t_j + 1) = t_i^{BR}(t_j)$  or  $t_i^{BR}(t_j + 1) = t_i^{BR}(t_j) - 1$ . Suppose that  $(k, l)$  constitutes a reactive SSDE. Given the behavior of  $t_D^{BR}(t_H)$ , only the four following threshold profiles may also constitute reactive SSDEs:  $(k - 1, l + 1), (k - 1, l), (k + 1, l)$  or to  $(k + 1, l - 1)$ . Furthermore, given the behavior of  $t_H^{BR}(t_D)$ , only the four following threshold profiles may also constitute reactive SSDEs:  $(k - 1, l + 1), (k, l + 1), (k, l - 1)$  or  $(k + 1, l - 1)$ . Taking



the intersection of the two sets, the only neighbouring points to  $(k, l)$  that may constitute reactive SSDEs are  $(k - 1, l + 1)$  or  $(k + 1, l - 1)$ . Suppose finally that the two best response functions do not intersect in any of these two neighbouring points. Then, this implies that they do not intersect in any other point than  $(k, l)$ . ■

### B.3 Normative Analysis

#### Proof of Proposition 2.1 (reactive SPDE vs reactive SNDE):

**Step 1** Recall that the unique reactive SPDE involves doves truthfully revealing their signal and voting according to  $T_D^{n_D}$  while hawks babble and always convict.

**Step 2** Recall that there always exists a unique reactive SNDE, given by profile A or B. Recall also that profile A is given by  $(\sigma_g^H = 1, \sigma_i^H = 1, \sigma_g^D = 1, \sigma_i^D = 0)$ . Suppose that  $\beta(p, n_D - 1, n_D) < q_D < \beta(p, n_D, n_D)$ , so that equilibrium A is the unique reactive SNDE. For these parameter values, the unique reactive SNDE and the unique reactive SPDE are thus outcome equivalent.

**Step 3** Steps 3 to 9 are dedicated to the examination of parameter values for which profile B is the unique reactive SNDE (i.e. iff  $1 - p < q_D < \beta(p, n_D - 1, n_D)$ ). Recall that the latter equilibrium is given by  $(\sigma_g^H = 1, \sigma_i^H = 1, \sigma_g^D = 1, \sigma_i^D = y)$ , with  $y \in (0, 1)$ . The unique reactive SPDE is here characterized by a dove threshold  $T_D^{n_D} \leq n_D - 1$ . The transition from the unique reactive SNDE to the unique SPDE is clearly strictly beneficial to the doves, as these are now optimally aggregating their information. In contrast, it however remains unclear whether the transition from the first to the second equilibrium is strictly beneficial to the hawks as well. If we can prove that this is the case, then we know that the unique reactive SPDE is strongly Pareto improving w.r.t to the unique reactive SNDE, for the concerned parameter values.

**Step 3** All we need is thus to show that, starting from the reactive SND profile B, allowing doves to Subgroup Deliberate while keeping the hawks' play fixed will be strictly beneficial to the hawks. We do so in the next steps. Denote by  $\Pi_j(q_D, SD, t_D)$  the expected payoff of preference type  $j$  when the doves are allowed to Subgroup Deliberate and adopt a threshold  $t_D$ , while hawks always all vote for conviction as in the reactive SND profile B. Let  $t_D(q_D)$  be the optimal threshold adopted by the doves in these circumstances, given  $q_D$ , i.e. let  $t_D(q_D) = T_D^{n_D}$ . Denote by  $\Pi_j(q_D, ND)$  the expected payoff of preference type  $j$  in the

reactive SND equilibrium B. Denote by  $y(q_D)$  the mixing probability of the doves after an  $i$ -signal in the reactive SND equilibrium B. Note that:

$$\begin{aligned}
W(q_j, q_D) &\equiv \Pi_j(q_D, SD, t_D(q_D)) - \Pi_j(q_D, ND) & (B.24) \\
&= -P(G) \sum_{x=0}^{n_D} B(p, x, n_D) [y(q_D)]^{n_D-x} (1 - q_j) \\
&\quad + P(I) \sum_{x=0}^{n_D} B(1 - p, x, n_D) [y(q_D)]^{n_D-x} q_j \\
&\quad + P(G) \sum_{x=t_D(q_D)}^{n_D} B(p, x, n_D) (1 - q_j) - P(I) \sum_{x=t_D(q_D)}^{n_D} B(1 - p, x, n_D) q_j. & (B.25)
\end{aligned}$$

It follows that:

$$\begin{aligned}
\partial W(q_j, q_D) / \partial q_j &= \frac{1}{2} \sum_{x=0}^{n_D} (B(p, x, n_D) + B(1 - p, x, n_D)) [y(q_D)]^{n_D-x} & (B.26) \\
&\quad - \frac{1}{2} \sum_{x=t_D(q_D)}^{n_D} (B(p, x, n_D) + B(1 - p, x, n_D)).
\end{aligned}$$

The sign of  $\partial W(q_j, q_D) / \partial q_j$  is thus determined by the difference in the total probability of conviction implied by each of the two voting scenarios considered, i.e. No Deliberation by the doves according to the symmetric voting strategy ( $\sigma_g^D = 1, \sigma_i^D = y(q_D)$ ) or Subgroup Deliberation by the doves with an optimally chosen conviction threshold  $t_D(q_D)$ . As the hawks' strategy is unchanged and the doves are able to share their information when they Subgroup Deliberate,  $W(q_D, q_D) > 0$ . If we can show that for all values of  $q_D$  and corresponding values  $t_D(q_D)$  and  $y(q_D)$ , the derivative  $\partial W(q_j, q_D) / \partial q_j$  is negative, then it is also true that  $W(q_H, q_D) > 0$ , because  $q_H < q_D$ . Which in other words means that also the hawks benefit from the change in the doves' strategy, if they continue to apply the strategy ( $\sigma_g^H = 1, \sigma_i^H = 1$ ) that they follow in the reactive SND equilibrium B.

**Step 4** Define the following two expressions:

$$I(n_D) \equiv \left\{ \frac{n_D}{2} + 1 \text{ if } n_D \text{ is even, } \frac{n_D + 1}{2} \text{ if } n_D \text{ is uneven.} \right\} \quad (B.27)$$

and for all  $z \in \{I(n_D), \dots, n_D\}$

$$\Psi(z) \equiv \begin{cases} \partial W(q_j, \frac{1}{2}) / \partial q_j & \text{for } z = I(n_D) \text{ and } n_D \text{ uneven,} \\ \lim_{\varepsilon \rightarrow 0^+} \partial W(q_j, \beta(p, z-1, n_D) + \varepsilon) / \partial q_j & \text{otherwise.} \end{cases} \quad (\text{B.28})$$

In order to show that  $\partial W(q_j, q_D) / \partial q_j$  is negative for all  $q_D \in (\frac{1}{2}, \beta(p, n_D - 1, n_D))$ , it is enough to verify that  $\Psi(z) \leq 0$ , for all  $z \in \{I(n_D), \dots, n_D\}$ . This is true for the two following reasons. First, stating that  $\Psi(z) \leq 0$ , for all  $z \in \{I(n_D), \dots, n_D\}$  is equivalent to stating that  $\partial W(q_j, q_D) / \partial q_j \leq 0$  for  $q_D = \frac{1}{2}$  as well as for  $q_D = \lim_{\varepsilon \rightarrow 0^+} \beta(p, z-1, n_D) + \varepsilon$ ,  $\forall z \in \{I(n_D) + 1, \dots, n_D\}$ . Secondly, given that  $y(q_D)$  is decreasing in  $q_D$  and given that  $t_D(q_D)$  is constant for all  $q_D \in (\beta(p, z-1, n_D), \beta(p, z, n_D)]$ , the derivative  $\partial W(q_j, q_D) / \partial q_j$  is a decreasing function of  $q_D$  for all  $q_D \in (\beta(p, z-1, n_D), \beta(p, z, n_D)]$ .

**Step 5** The proof that  $\Psi(z) \leq 0$  for all  $z \in \{I(n_D), \dots, n_D\}$  is divided into five steps (6, 7, 8, 9 and 10). Step 6 shows that  $\Psi(n_D) \leq 0$ . Step 7 shows that  $\Psi(I(n_D)) \leq 0$ , for all  $n_D$  even. Step 8 shows that  $\Psi(I(n_D)) \leq 0$  and  $\Psi(I(n_D) + 1) \leq 0$ , for all  $n_D$  uneven. Step 8 shows the following. If  $n_D$  is even, then if  $\Psi(z) \leq \Psi(z+1)$ , it follows that  $\Psi(z+1) \leq \Psi(z+2)$  for all  $z \in \{I(n_D), \dots, n_D - 1\}$ . If, in contrast,  $n_D$  is uneven, then if  $\Psi(z) \leq \Psi(z+1)$ , it follows that  $\Psi(z+1) \leq \Psi(z+2)$  for all  $z \in \{I(n_D) + 1, \dots, n_D - 1\}$ . Step 10, finally, shows that the four facts proven in steps 6, 7, 8 and 9 imply together that  $\Psi(z) \leq 0$ , for all  $z \in \{I(n_D), \dots, n_D\}$ .

**Step 6** Note the following fact:

**Fact 1:**  $\Psi(n_D) < 0$  whether  $n_D$  is even or uneven.

Setting  $z = n_D$ , Fact 1 follows immediately from the fact that  $y(\beta(p, n_D - 1, n_D)) = 0$  while  $\lim_{\varepsilon \rightarrow 0^+} t_D(\beta(p, n_D - 1, n_D) + \varepsilon) = n_D$ .

**Step 7** Note the following fact:

**Fact 2:**  $\Psi(I(n_D)) < 0$  if  $n_D$  is even.

Note here that  $\beta(p, I(n_D) - 1, n_D) = \frac{1}{2}$ . Also,  $t_D(q_D) = I(n_D)$  if  $q_D \in (\frac{1}{2}, \beta(p, I(n_D), n_D))$ . For  $t_D(q_D) = I(n_D)$ , the total probability of conviction, if doves Subgroup Deliberate and

hawks always convict, is given by:

$$\frac{1}{2} \sum_{x=I(n_D)}^{n_D} (B(p, x, n_D) + B(1-p, x, n_D)) = \frac{1}{2} (1 - B(p, \frac{n_D}{2}, n_D)). \quad (\text{B.29})$$

On the other hand, for  $q_D = \frac{1}{2}$ , the total probability of conviction in the equilibrium B is given by:

$$\frac{1}{2} \frac{(2p-1)^{n_D} \left( \left( \frac{(1-p)}{p} \right)^{\frac{n_D}{n_D-1}} + 1 \right)}{\left( p - \left( \frac{(1-p)}{p} \right)^{\frac{1}{n_D-1}} (1-p) \right)^{n_D}}. \quad (\text{B.30})$$

Now, note that  $(\text{B.30}) \leq (\text{B.29})$ , for any  $p > \frac{1}{2}$  and  $n_D \geq 4$ . Note that given that we impose  $q_D > \frac{1}{2}$ , the equilibrium B does not exist if  $n_D = 2$  so that we can ignore this case. Indeed, B exists only if  $q_D < \beta(p, n_D - 1, n_D)$ . For the case of  $n_D = 2$ , this translates into  $q_D < \beta(p, 1, 2) = \frac{1}{2}$  which contradicts the assumption that  $q_D > \frac{1}{2}$ .

**Step 8** Note the following fact:

**Fact 3:**  $\Psi(I(n_D)) < 0$  and  $\Psi(I(n_D) + 1) < 0$  if  $n_D$  is uneven.

We first look at  $\Psi(I(n_D))$ . For  $q_D = \frac{1}{2}$  note that  $t_D(q_D) = I(n_D)$ . The total probability of conviction for  $t_D(q_D) = I(n_D)$ , if doves Subgroup Deliberate and hawks always convict, is given by:

$$\frac{1}{2} \sum_{x=I(n_D)}^{n_D} (B(p, x, n_D) + B(1-p, x, n_D)) = \frac{1}{2}. \quad (\text{B.31})$$

On the other hand, for  $q_D = \frac{1}{2}$ , the total probability of conviction in the equilibrium B is given by:

$$\frac{1}{2} \frac{(2p-1)^{n_D} \left( \left( \frac{(1-p)}{p} \right)^{\frac{n_D}{n_D-1}} + 1 \right)}{\left( p - \left( \frac{(1-p)}{p} \right)^{\frac{1}{n_D-1}} (1-p) \right)^{n_D}}. \quad (\text{B.32})$$

We now look at  $\Psi(I(n_D) + 1)$ . Note that  $t_D(q_D) = I(n_D) + 1$  if

$$q_D \in (\beta(p, I(n_D), n_D), \beta(p, I(n_D) + 1, n_D)).$$

The total probability of conviction for  $t_D(q_D) = I(n_D) + 1$ , if doves Subgroup Deliberate and

hawks always convict, is given as follows :

$$\frac{1}{2} \sum_{x=I(n_D)+1}^{n_D} (B(p, x, n_D) + B(1-p, x, n_D)) = \frac{1}{2} (1 - B(p, I(n_D), n_D) - B(1-p, I(n_D), n_D)). \quad (\text{B.33})$$

On the other hand, for  $q_D = \beta(p, I(n_D), n_D)$ , the total probability of conviction in the equilibrium B is given by:

$$\frac{1}{2} \frac{(2p-1)^{n_D} \left( \left( \frac{(1-p)^2}{p^2} \right)^{\frac{n_D}{n_D-1}} + 1 \right)}{\left( p - \left( \frac{(1-p)^2}{p^2} \right)^{\frac{1}{n_D-1}} (1-p) \right)^{n_D}}. \quad (\text{B.34})$$

Now, note that  $(B.32) < (B.31)$  and  $(B.34) \leq (B.33)$ , for any  $p \in (\frac{1}{2}, 1]$  and  $n_D \geq 3$ . Note that for  $n_D = 1$ , the equilibrium B does not exist so that this case can be ignored. Indeed, B exists only if  $q_D \leq \beta(p, 0, 1) = 1-p$  if  $n_D = 1$ . But we have assumed  $q_D > \frac{1}{2}$ .

**Step 9** Note the following fact:

**Fact 4:** If  $\Psi(z+1) - \Psi(z) > 0$  then  $\Psi(z+2) - \Psi(z+1) > 0$ ,  
for all  $z \in \{I(n_D), \dots, n_D - 1\}$  if  $n_D$  even,  
for all  $z \in \{I(n_D) + 1, \dots, n_D - 1\}$  if  $n_D$  uneven.

Using the Binomial Formula, for  $q_D = \beta(p, z-1, n_D)$ , we may define and rewrite the following new function, which we use to prove the statement:

$$\begin{aligned} \Phi(p, z, n_D) &\equiv \left( \sum_{x=0}^{n_D} (B(p, x, n_D) + B(1-p, x, n_D)) \right) [y(\beta(p, z-1, n_D))]^{n_D-x} \\ &= \frac{(2p-1)^{n_D} \left( \left( \frac{B(1-p, z-1, n_D)(1-p)}{B(p, z-1, n_D)p} \right)^{\frac{n_D}{n_D-1}} + 1 \right)}{\left( p - \left( \frac{B(1-p, z-1, n_D)(1-p)}{B(p, z-1, n_D)p} \right)^{\frac{1}{n_D-1}} (1-p) \right)^{n_D}}. \end{aligned} \quad (\text{B.35})$$

Note that:

$$\begin{aligned} \Psi(z+1) - \Psi(z) &= \Phi(p, z+1, n_D) - \Phi(p, z, n_D) \\ &\quad + B(p, z-1, n_D) + B(1-p, z-1, n_D). \end{aligned} \quad (\text{B.36})$$

Also,

$$B(p, z - 1, n_D) + B(1 - p, z - 1, n_D) > 0, \quad \forall z \in \{1, \dots, n_D\}. \quad (\text{B.37})$$

Note furthermore that

$$\frac{1}{2}\Phi(p, z, n_D) + \frac{1}{2}\Phi(p, z + 2, n_D) > \Phi(p, z + 1, n_D). \quad (\text{B.38})$$

Inequality (B.38) follows from the fact that the function  $\Phi(p, z, n_D)$  is decreasing and convex in  $z$  over the relevant domain. The latter fact follows from the fact that the following two functions:

$$f_1(p, n_D, z) \equiv \left( \frac{B(1 - p, z - 1, n_D)(1 - p)}{B(p, z - 1, n_D)p} \right)^{\frac{n_D}{n_D - 1}} + 1 \quad (\text{B.39})$$

and

$$f_2(p, n_D, z) \equiv \frac{1}{\left( p - \left( \frac{B(1 - p, z - 1, n_D)(1 - p)}{B(p, z - 1, n_D)p} \right)^{\frac{1}{n_D - 1}} (1 - p) \right)^{n_D}}. \quad (\text{B.40})$$

are themselves decreasing and convex in  $z$  over the relevant domain. Note finally that:

$$\frac{1}{2}\Phi(p, z, n_D) + \frac{1}{2}\Phi(p, z + 2, n_D) > \Phi(p, z + 1, n_D) \quad (\text{B.41})$$

$$\Leftrightarrow \Phi(p, z + 1, n_D) - \Phi(p, z, n_D) < \Phi(p, z + 2, n_D) - \Phi(p, z + 1, n_D).$$

Using (B.36),(B.37),(B.38),(B.41) yields our statement that  $\Psi(z + 2) - \Psi(z + 1)$  is also positive whenever  $\Psi(z + 1) - \Psi(z)$  is positive.

**Step 10** From Facts 1,2 and 3 we know that  $\Psi(z)$  is negative at the boundaries. From Fact 4, we know that if  $\Psi(z)$  starts to increase it never decreases again. It follows that it has to be that  $\Psi(z) \leq 0$ , for all  $z \in \{I(n_D), \dots, n_D\}$ , whether  $n_D$  is even or uneven.

**Step 11** Given that  $\Psi(z) \leq 0$ , for all  $z \in \{I(n_D), \dots, n_D\}$ , it follows by the argument given in step 4 that  $\partial W(q_j, q_D)/\partial q_j \leq 0$  for all  $q_D \in (\frac{1}{2}, \beta(p, n_D - 1, n_D))$ , which implies that  $W(q_H, q_D) > 0$  for all  $q_H \in [0, q_D)$  and  $q_D \in (\frac{1}{2}, \beta(p, n_D - 1, n_D))$ . ■

**Proof of Proposition 2.2 (reactive SSDEs):**

This complements the part of the proof of Proposition 2.2 that appears in the main text. We prove in what follows that transiting from  $(t_H - 1, t_D + 1)$  to  $(t_H - 1, t_D)$  is beneficial

for the preference type  $H$  given our assumption that  $m > 1$ . A similar argument shows that transiting from  $(t_H - 1, t_D + 1)$  to  $(t_H, t_D + 1)$  is beneficial for the preference type  $D$  given our assumption that  $m > 1$ . Assume that

$$\frac{B(p, t_D, n_D) \left( \sum_{x \geq t_{H-1}}^{n_H} B(p, x, n_H) \right)}{B(1-p, t_D, n_D) \left( \sum_{x \geq t_{H-1}}^{n_H} B(1-p, x, n_H) \right)} < \frac{q_H}{1 - q_H} \quad (\text{B.42})$$

and

$$\frac{q_D}{1 - q_D} \leq \frac{B(p, t_D + 1, n_D) \left( \sum_{x \geq t_{H-1}}^{n_H} B(p, x, n_H) \right)}{B(1-p, t_D + 1, n_D) \left( \sum_{x \geq t_{H-1}}^{n_H} B(1-p, x, n_H) \right)}. \quad (\text{B.43})$$

By a standard argument already used in Appendix B.2, we furthermore know that by definition, there exists some integer  $T \in \{1, \dots, n\}$  s.t.

$$\frac{B(p, T - 1, n)}{B(1-p, T - 1, n)} < \frac{B(p, t_D + 1, n_D) \left( \sum_{x \geq t_{H-1}}^{n_H} B(p, x, n_H) \right)}{B(1-p, t_D + 1, n_D) \left( \sum_{x \geq t_{H-1}}^{n_H} B(1-p, x, n_H) \right)} \leq \frac{B(p, T, n)}{B(1-p, T, n)} \quad (\text{B.44})$$

and

$$\frac{B(p, T - 2, n)}{B(1-p, T - 2, n)} < \frac{B(p, t_D, n_D) \left( \sum_{x \geq t_{H-1}}^{n_H} B(p, x, n_H) \right)}{B(1-p, t_D, n_D) \left( \sum_{x \geq t_{H-1}}^{n_H} B(1-p, x, n_H) \right)} \leq \frac{B(p, T - 1, n)}{B(1-p, T - 1, n)}. \quad (\text{B.45})$$

Now, the inequalities (B.42), (B.43), (B.44) and (B.45) imply that there is some integer  $T \in \{1, \dots, n\}$  s.t.

$$\frac{B(p, T - 2, n)}{B(1-p, T - 2, n)} < \frac{q_H}{1 - q_H} < \frac{q_D}{1 - q_D} \leq \frac{B(p, T, n)}{B(1-p, T, n)},$$

which contradicts our assumption that  $m > 1$ . It follows that (B.42) and (B.43) cannot be true. ■

### **Proof of Proposition 2.3 (reactive SSDEs vs reactive SPDE):**

**Step 1** The unique reactive SPDE is characterized by a dove threshold  $T_D^{nD}$ . Now, there are two cases to analyze (a. and b.).

In Case a),  $q_H \leq P(G | |g|_H = 0, |g|_D \geq T_D^{nD})$  and there exists a reactive SSDE given by  $t_H = 0$  and  $t_D = T_D^{nD}$ . This latter reactive SSDE is outcome equivalent to the unique reactive

simple SPDE. If there exists any other reactive SSDE, then by Proposition 2.2, it is strongly Pareto dominant w.r.t. the reactive SSDE in which  $t_H = 0$  and  $t_D = T_D^{n_D}$ , and thus also strongly Pareto dominant w.r.t. the unique reactive SPDE.

**Step 2** In Case b),  $q_H > P(G | |g|_H = 0, |g|_D \geq T_D^{n_D})$  and there thus exists no reactive SSDE given by  $t_H = 0$  and  $t_D = T_D^{n_D}$ . We know however from Lemma 2.6 that there exists some reactive SSDE. We now conduct an argument based on a hypothetical adjustment process. Start from the reactive SSD profile in which  $t_H = 0$  and  $t_D = T_D^{n_D}$ . We know that this profile (although it is not an equilibrium profile) yields a payoff to each preference type that is equivalent to that received in the unique reactive SPDE. Now, let hawks choose their collective best response to  $T_D^{n_D}$ , i.e.  $t_H^{BR}(T_D^{n_D})$ . We know that the latter is strictly larger than 0 given that  $q_H > P(G | |g|_H = 0, |g|_D \geq T_D^{n_D})$ . This adjustment is strictly beneficial to hawks and also to doves, given that hawks become more lenient. In a further step, let doves revise their threshold and choose their own best response  $t_D^{BR}(t_H^{BR}(T_D^{n_D}))$ . Again, the adjustment is by definition beneficial to doves as well as to hawks, as doves become weakly harsher. Repeat the adjustment of the hawks, etc.

This process of mutual adjustment converges to a reactive SSDE, and every step of the adjustment process is strictly welfare improving for both preference types. It follows that the reactive SSDE to which our adjustment process converges is strongly Pareto dominant w.r.t. the unique reactive SPDE. Note furthermore that any other reactive SSDE is less polarized than this first reactive SSDE and thus, by Proposition 2.2, strongly Pareto improving w.r.t. the latter. It follows that any reactive SSDE is strongly Pareto dominant w.r.t. the unique reactive SPDE. ■







## Appendix C

# Carbon Taxation under Asymmetric Information over Fossil-fuel Reserves

We solve most maximization problems by the first order condition, which is valid as we deal with concave functions. We do not always highlight when we use this method.

### C.1 The Case of Symmetric Information

The proof of Lemma 3.1 is postponed to after the proof of Proposition 3.1.

**Proof of Proposition 3.1:** We solve the symmetric information problem for any size of the reserves,  $Q \geq 0$ , and by backward induction.<sup>1</sup> Proposition 3.1 follows directly from these results.

**Without the quantity constraint:** Let us first solve the problem without the quantity constraint (3.13).

Given  $Q_1$  and  $\tau$  the fossil-fuel owner  $Q$ 's problem in period 2 is

$$Q_2^{**}(Q_1, \tau, Q) \equiv \operatorname{argmax}_{Q_2} U_F^{t=2}(Q_2, \tau).$$

---

<sup>1</sup>Let  $Q_2^{x**}(Q_1, \tau) \equiv Q_2^{**}(Q_1, \tau, Q^x)$ ,  $\tau^{x**}(Q_1) \equiv \tau^{**}(Q_1, Q^x)$  and  $Q_1^{x**} = Q_1^{**}(Q^x)$  for all  $x \in \{L, H\}$

The maximizer of the utility function  $U_F^{t=2}(Q_2) = (1 - Q_2 - \tau)Q_2$  is

$$Q_2^{**}(Q_1, \tau, Q) = \frac{1 - \tau}{2}.$$

The regulator's problem between period 1 and period 2, given,  $Q_1$ , and the best response,  $Q_2^{**}(\tau, Q_1, Q)$ , is hence

$$\tau^{**}(Q_1, Q) \equiv \operatorname{argmax}_{\tau} U_S^{t=2}(Q_1, \frac{1 - \tau}{2}, \tau).$$

Rewriting, we solve

$$\begin{aligned} \tau^{**}(Q_1, Q) &= \operatorname{argmax}_{\tau} \left[ \lambda \left(1 - \frac{1 - \tau}{2}\right) \left(\frac{1 - \tau}{2}\right) + \frac{1}{2}(1 - \lambda) \left(\frac{1 - \tau}{2}\right)^2 \right. \\ &\quad \left. + (1 - 2\lambda) \left(\frac{1 - \tau}{2}\right) \tau - \frac{1 + r}{r} \frac{d}{2} (Q_1 + \frac{1 - \tau}{2})^2 \right] \\ &= 1 + \frac{2dQ_1(1 + r) - 2(1 - \lambda)r}{3r - 5\lambda r + d(1 + r)} \equiv \tau_{op}(Q_1). \end{aligned}$$

Taking the results together the fossil-fuel owner's problem in period 1 is hence

$$Q_1^{**}(Q) \equiv \operatorname{argmax}_{Q_1} U_F(Q_1, \frac{1 - \tau_{op}(Q_1)}{2}, \tau_{op}(Q_1)).$$

By the first order condition we find

$$Q_1^{**}(Q) = \frac{1}{2} \cdot \frac{(3r - 5\lambda r + d(1 + r))^2 - 2d(1 - \lambda)r}{(3r - 5\lambda r + d(1 + r))^2 - d^2(1 + r)} \equiv Q_1^{op}.$$

As we have solved the problem without the quantity constraint we can now deduce that whenever the reserves,  $Q$ , are high enough

$$Q \geq Q_k^1 \equiv Q_k^1(d, r, \lambda) \equiv Q_1^{op} + \frac{1 - \tau_{op}(Q_1^{op})}{2}$$

the equilibrium strategies are

$$Q_1^{**}(Q) = Q_1^{op}, \quad \tau^{**}(Q) = \tau_{op}(Q_1^{op}), \quad Q_2^{**}(Q) = \frac{1 - \tau_{op}(Q_1^{op})}{2}.$$

**With the quantity constraint:** Now assume  $Q < Q_k^1$ . The fossil-fuel owner is not able

to supply the above given amounts as  $\frac{1-\tau_{op}(Q_1^{op})}{2} > Q - Q_1^{op}$ , i.e. the quantity constraint is binding,  $Q_2^{**}(Q_1^{op}, \tau_{op}(Q_1^{op}), Q) = Q - Q_1^{op}$ . Let us now solve the problem with the quantity constraint (3.13).

**Period 2:** Given  $Q_1$  and  $\tau$  the fossil-fuel owner  $Q$ 's problem in period 2 is

$$Q_2^{**}(Q_1, \tau, Q) \equiv \operatorname{argmax}_{Q_2} U_F^{t=2}(Q_2, \tau)$$

s.t.

$$Q_2^{**}(Q_1, \tau, Q) \leq Q - Q_1.$$

As the maximizer of the utility function  $U_F^{t=2}(Q_2)$  is  $Q_2 = \frac{1-\tau}{2}$  the solution to the above problem is simply

$$Q_2^{**}(Q_1, \tau, Q) = \min\left\{\frac{1-\tau}{2}, Q - Q_1\right\} = \begin{cases} \frac{1-\tau}{2} & \text{if } \tau \geq \tau_c(Q_1, Q) \\ Q - Q_1 & \text{if } \tau \leq \tau_c(Q_1, Q) \end{cases},$$

where  $\tau_c(Q_1, Q) \equiv 1 - 2Q + 2Q_1$ .<sup>2</sup>

**Between period 1 and period 2:** The regulator's problem given,  $Q_1$ , and the best response,  $Q_2^{**}(\tau, Q_1, Q)$ , is therefore

$$\begin{aligned} \tau^{**}(Q_1, Q) &\equiv \operatorname{argmax}_{\tau} U_S^{t=2}(Q_1, Q_2^{**}(Q_1, \tau, Q), \tau) \\ &= \operatorname{argmax}_{\tau} \begin{cases} U_S^{t=2}(Q_1, \frac{1-\tau}{2}, \tau) & \text{if } \tau \geq \tau_c(Q_1, Q) \\ U_S^{t=2}(Q_1, Q - Q_1, \tau) & \text{if } \tau \leq \tau_c(Q_1, Q) \end{cases}. \end{aligned}$$

Before we derive the solution assume  $Q_2^{**}(Q_1, \tau, Q) = Q - Q_1$ . The regulator's problem transforms to

$$\begin{aligned} \tau^{**}(Q_1, Q) &= \operatorname{argmax}_{\tau} \left[ \lambda(1 - Q + Q_1)(Q - Q_1) + \frac{1}{2}(1 - \lambda)(Q - Q_1)^2 \right. \\ &\quad \left. + (1 - 2\lambda)(Q - Q_1)\tau - \frac{1+r}{r} \frac{d}{2} Q^2 \right] \equiv \tau_b(Q_1, Q). \end{aligned} \tag{C.1}$$

---

<sup>2</sup>Note that  $Q_2^{**}(\tau, Q_1, Q)$  is a decreasing function of the tax  $\tau$ .

Note that the argument

$$\lambda(1 - Q + Q_1)(Q - Q_1) + \frac{1}{2}(1 - \lambda)(Q - Q_1)^2 + (1 - 2\lambda)(Q - Q_1)\tau - \frac{1 + r}{r} \frac{d}{2} Q^2$$

is a linear and increasing function in  $\tau$ . Recall that the fossil-fuel owner exhausts his reserves,  $Q_2^{**}(Q_1, \tau, Q) = Q - Q_1$ , if and only if the tax,  $\tau$ , satisfies  $\tau \leq \tau_c(Q_1, Q)$ ; and he does not exhaust his reserves,  $Q_2^{**}(Q_1, \tau, Q) = \frac{1-\tau}{2}$ , if the tax instead satisfies  $\tau > \tau_c(Q_1, Q)$ . So, as the above problem is increasing in  $\tau$  the regulator best response is  $\tau_b(Q_1, Q) = \tau_c(Q_1, Q)$ .

We now solve the regulator's general problem. Recall again that for any supply,  $Q_1$ , in period 1 the regulator forces the fossil-fuel owner to exhaust his reserves,  $Q - Q_1$ , in period 2 if she sets a tax  $0 \leq \tau \leq \tau_c(Q_1, Q)$ , and she forces the fossil-fuel owner not to exhaust his reserves,  $\frac{1-\tau}{2}$ , if she instead sets a tax  $\tau > \tau_c(Q_1, Q)$ . We divide the following analysis into two cases.

1. If  $Q_1$  is such that  $\tau_{op}(Q_1) > \tau_c(Q_1, Q)$ , the fossil-fuel owner's response given the tax,  $\tau_{op}(Q_1)$ , satisfies  $Q_2^{**}(Q_1, \tau_{op}(Q_1), Q) = \frac{1-\tau_{op}(Q_1)}{2}$ . First, as  $\tau_{op}(Q_1)$  is the maximizer of the regulator's tax problem without the quantity constraint,  $\tau_{op}(Q_1) = \operatorname{argmax}_{0 \leq \tau \leq 1} U_S^{t=2}(Q_1, \frac{1-\tau}{2}, \tau)$ , second as  $\tau_{op}(Q_1) > \tau_c(Q_1, Q)$ , and third as  $U_S^{t=2}(Q_1, \frac{1-\tau}{2}, \tau) \geq U_S^{t=2}(Q_1, Q - Q_1, \tau)$  for all  $\tau \in [0, 1]$ , we can straightforwardly conclude that the regulator's optimal response to the supply,  $Q_1$ , is

$$\tau^{**}(Q_1, Q) = \tau_{op}(Q_1) \quad \text{if } \tau_{op}(Q_1) > \tau_c(Q_1, Q).$$

2. If  $Q_1$  is such that  $\tau_{op}(Q_1) \leq \tau_c(Q_1, Q)$ , we show that the regulator's utility is maximized at  $\tau_c(Q_1, Q)$ . Indeed, the fossil-fuel owner's response given the tax,  $\tau_{op}(Q_1)$ , satisfies  $Q_2^{**}(Q_1, \tau_{op}(Q_1), Q) = Q - Q_1$ . As the regulator prefers tax revenues whenever the reserves are exhausted (see the above analysis, equation (C.1)) she prefers to set  $\tau_c(Q_1, Q)$  instead of any  $0 \leq \tau \leq \tau_c(Q_1, Q)$  and so also instead of  $\tau_{op}(Q_1)$ , i.e.  $\tau_c(Q_1, Q) = \operatorname{argmax}_{\tau \leq \tau_c(Q_1, Q)} U_S^{t=2}(Q_1, Q - Q_1, \tau)$ . Now, what about a tax  $\tau > \tau_c(Q_1, Q)$ . As the regulator's utility,  $U_S^{t=2}(Q_1, \frac{1-\tau}{2}, \tau)$ , is concave in  $\tau$  and maximized at  $\tau_{op}(Q_1)$  and  $\tau_{op}(Q_1) \leq \tau_c(Q_1, Q)$  she prefers to set the tax at the boundary,  $\tau_c(Q_1, Q)$ , followed by the fossil-fuel owner's best response,  $Q_2^{**}(Q_1, \tau_c(Q_1, Q), Q) = \frac{1-\tau_c(Q_1, Q)}{2} = Q - Q_1$ , to any tax  $\tau > \tau_c(Q_1, Q)$  followed by the best response,  $Q_2^{**}(Q_1, \tau, Q) = \frac{1-\tau}{2}$ , i.e.

$\tau_c(Q_1, Q) = \operatorname{argmax}_{\tau_c(Q_1, Q) \leq \tau} U_S^{t=2}(Q_1, \frac{1-\tau}{2}, \tau)$ . So, the regulator's optimal response to  $Q_1$  is

$$\tau^{**}(Q_1, Q) = \tau_c(Q_1, Q) \text{ if } \tau_{op}(Q_1) \leq \tau_c(Q_1, Q).$$

Summarizing we can state that the regulator's best tax response satisfies

$$\tau^{**}(Q_1, Q) = \max \{ \tau_{op}(Q_1), \tau_c(Q_1, Q) \} = \begin{cases} \tau_c(Q_1, Q) & \text{if } Q_1 \geq Q_1^B(Q) \\ \tau_{op}(Q_1) & \text{if } Q_1 \leq Q_1^B(Q) \end{cases},$$

$$\text{with } Q_1^B(Q) \equiv \frac{3Qr - 5Q\lambda r + Qd(1+r) - (1-\lambda)r}{3r - 5\lambda r}.$$

**Period 1:** The fossil-fuel owner's problem in period 1 is to supply an amount,  $Q_1$ , which maximizes his profits given the expected responses of all players in the following periods,  $\tau^{**}(Q_1, Q)$  and  $Q_2^{**}(Q_1, \tau(Q_1, Q), Q)$ . In other words

$$\begin{aligned} Q_1^{**}(Q) &= \operatorname{argmax}_{Q_1} U_F(Q_1, Q_2^{**}(Q_1, \tau^{**}(Q_1, Q), Q), \tau^{**}(Q_1, Q)) \\ &= \operatorname{argmax}_{Q_1} \begin{cases} U_F(Q_1, Q - Q_1, \tau_c(Q_1, Q)) & \text{if } Q_1 \geq Q_1^B(Q) \\ U_F(Q_1, \frac{1-\tau_{op}(Q_1)}{2}, \tau_{op}(Q_1)) & \text{if } Q_1 \leq Q_1^B(Q) \end{cases}. \end{aligned}$$

Note that for different strategies,  $Q_1$ , in period 1 the fossil-fuel owner exposes himself to two qualitatively different situations. Either given the tax response of the social planner,  $\tau^{**}(Q_1, Q)$ , he exhausts his reserves,  $Q_2^{**}(Q_1, \tau, Q) = Q - Q_1$ , or he does not exhaust his reserves  $Q_2^{**}(Q_1, \tau, Q) = \frac{1-\tau}{2}$ .

As  $Q \leq Q_k^1$  we have  $Q_1^{op} \geq Q_1^B(Q)$ . Together with the fact that  $Q_1^{op}$  maximizes the regulator's utility,  $U_F(Q_1, \frac{1-\tau_{op}(Q_1)}{2}, \tau_{op}(Q_1))$ , and as  $Q_1^{op} \geq Q_1^B(Q)$  we can infer that

$$U_F(Q_1^B(Q), \frac{1-\tau_{op}(Q_1^B(Q))}{2}, \tau_{op}(Q_1^B(Q))) = \max_{Q_1 \leq Q_1^B(Q)} U_F(Q_1, \frac{1-\tau_{op}(Q_1)}{2}, \tau_{op}(Q_1)).$$

Before we proceed let us solve

$$Q_1(Q) = \operatorname{argmax}_{Q_1} U_F(Q_1, Q - Q_1, \tau_b(Q_1, Q)).$$

The first order condition yields

$$Q_1(Q) = \frac{(1+r) - 2Q}{2r} \equiv Q_1^b(Q).$$

Taking the last two results together we can infer that  $Q_1^b(Q)$  maximizes the regulator's welfare at period 1 if and only if  $Q_1^B(Q) \leq Q_1^b(Q) \leq Q$ . So, if first  $Q_1^b(Q)$  can be supplied (the reserves are not allowed to be too low,  $Q \geq Q_1^b(Q)$ ) and if second given the supply,  $Q_1^b(Q)$ , in period 1 the regulator indeed wants the reserves to be exhausted (the reserves are not allowed to be too high, i.e. so that the tax indeed satisfies  $\tau^{**}(Q_1^b, Q) = \tau_b(Q_1^b(Q), Q)$  and consequently the response in period 2 is  $Q_2^{**}(Q_1^b(Q), \tau_b(Q_1^b(Q), Q), Q) = Q - Q_1^b(Q)$ ).

The first condition,  $Q_1^b(Q) \leq Q$ , is equivalent to  $Q \geq \frac{1}{2}$ . The second condition is equivalent to  $Q_1^B(Q) \leq Q_1^b(Q)$  which is then again the same as  $Q \leq \frac{(1+r)(3-5\lambda)+2(1-\lambda)r}{2(1+r)(3-5\lambda+d)} \equiv Q_k^{2,3}$ .

We can now conclude that if  $Q \leq Q_k^1$  the equilibrium strategies satisfy  $Q_2^{**}(Q) = Q - Q_1^{**}(Q)$ ,  $\tau^{**}(Q) = 1 + 2Q_1^{**}(Q) - 2Q$  and

1. if  $Q \leq \frac{1}{2}$

$$Q_1^{**}(Q) = Q, \quad Q_2^{**}(Q) = 0 \quad \text{and} \quad \tau^{**}(Q) = 1,$$

2. if  $\frac{1}{2} \leq Q \leq Q_k^2$

$$Q_1^{**}(Q) = \frac{(1+r) - 2Q}{2r}, \quad Q_2^{**}(Q) = \frac{(1+r)(2Q - 1)}{2r},$$

3. if  $Q_k^2 \leq Q$

$$Q_1^{**}(Q) = \frac{3Qr - 5Q\lambda r + Qd(1+r) - (1-\lambda)r}{3r - 5\lambda r}, \quad Q_2^{**}(Q) = \frac{(1-\lambda)r - Qd(1+r)}{3r - 5\lambda r}.$$

The comparative statics results and all other statements follow directly from these results. ■

**Proof of Lemma 3.1:** We show this result in two steps. Recall that  $\tau^{**}(Q_1, Q) = \max\{\tau_{op}(Q_1), \tau_b(Q_1, Q)\}$ :

1. Assume  $Q^L$  satisfies  $\tau^{L**}(Q_1) = \tau_{op}(Q_1)$ : As  $\tau^{L**}(Q_1) = \max\{\tau_{op}(Q_1), \tau_b(Q_1, Q^L)\}$  we can conclude that  $\tau_b(Q_1, Q^L) \leq \tau_{op}(Q_1)$ . In addition the inequality  $\tau_b(Q_1, Q^H) = 1 - 2Q^H + 2Q_1 < 1 - 2Q^L + 2Q_1 = \tau_b(Q_1, Q^L)$  holds. So we can furthermore conclude

---

<sup>3</sup>If  $0 < d < \frac{2(1-\lambda)r}{1+r}$  the thresholds are ordered as follows  $\frac{1}{2} \leq Q_k^2 \leq Q_k^1$ .



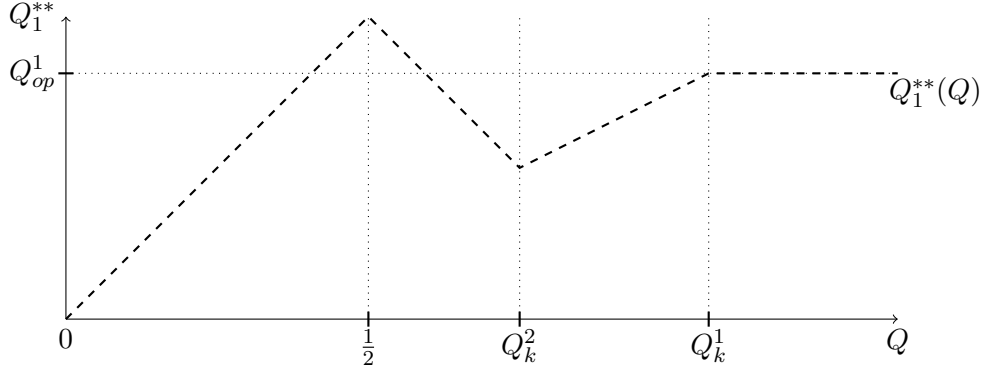


Figure C.1: Symmetric information equilibrium supply of period 1 as a function of the size of the reserves.

that  $\tau_b(Q_1, Q^H) \leq \tau_{op}(Q_1)$ , and hence

$$\tau^{H**}(Q_1) = \max\{\tau_{op}(Q_1), \tau_b(Q_1, Q^H)\} = \tau_{op}(Q_1) = \tau^{L**}(Q_1).$$

2. Assume  $Q^L$  satisfies  $\tau^{L**}(Q_1) = \tau_b(Q_1, Q^L)$ : We can conclude that  $\tau_{op}(Q_1) \leq \tau_b(Q_1, Q^L)$ :

(a) If  $Q^H$  satisfies  $\tau^{H**}(Q_1) = \tau_{op}(Q_1)$  then

$$\tau^{H**}(Q_1) \leq \tau^{L**}(Q_1).$$

(b) If  $Q^H$  instead satisfies  $\tau^{H**}(Q_1) = \tau_b(Q_1, Q^H)$  and as  $\tau_b(Q_1, Q^H) < \tau_b(Q_1, Q^L)$  we can again conclude:

$$\tau^{H**}(Q_1) \leq \tau^{L**}(Q_1).$$

■

## C.2 The Case of Asymmetric Information

### Separating Equilibria:

**Proof of Proposition 3.2:** To the low type. Condition (3.23) for  $x = L$  is equivalent

to condition

$$U_{L,L}(Q_1^{L*}) \geq \begin{cases} U_{L,L}(Q_1^{L**}) & \text{if } Q_1^{H*} > Q^L, \\ \max\{U_{L,L}(Q_1^{L**}), U_{L,H}(Q_1^{H*})\} & \text{if } Q_1^{H*} \leq Q^L \end{cases},^4 \quad (\text{C.2})$$

as the low type cannot supply an amount  $Q_1 > Q^L$ .

This implies as  $Q_1^{L*} \neq Q_1^{H*}$  that condition (C.2) is equivalent to condition

$$Q_1^{L*} = Q_1^{L**}, \quad (\text{C.3})$$

plus,  $U_{L,L}(Q_1^{L**}) \geq U_{L,H}(Q_1^{H*})$ , if  $Q_1^{H*} \leq Q^L$ .

To the high type. First, assume  $Q_1^{H**} > Q^L$ .

Given the supply  $Q_1^{H**}$  in period 1 –on or off the equilibrium path– and given any belief system,  $\mu^* \in \Omega$ , the subsequent tax response is the H-Tax-Response,  $\tau^*(Q_1^{H**}) = \tau^{H**}$ . Additionally as  $U_{H,H}(Q_1^{H**}) = \max_{Q_1} U_{H,H}(Q_1)$ , and as the L-Tax-Response is higher than the H-Tax-Response, i.e.  $U_{H,L}(Q_1) \leq U_{H,H}(Q_1)$  for all  $0 \leq Q_1 \leq Q^H$ , it follows immediately that condition (3.23) for  $x = H$  is equivalent to condition

$$U_{H,H}(Q_1^{H*}) \geq U_{H,H}(Q_1^{H**}) \quad \text{if } Q_1^{H**} > Q^L,$$

i.e.

$$Q_1^{H*} = Q_1^{H**} \quad \text{if } Q_1^{H**} > Q^L. \quad (\text{C.4})$$

Second, assume  $Q_1^{H**} \leq Q^L$ .

Note, as the high type's utility  $U_{H,H}(Q_1)$  is concave in  $Q_1$  and maximized at  $Q_1^{H**} \leq Q^L$  the highest high type's utility, given that, i) the supply of period 1 is off the equilibrium path, and ii) the subsequent tax is the H-Tax-Response, is his utility on the boundary,  $U_{H,H}(Q^L)$ . Also recall that  $U_{H,L}(Q_1^{L**}) = \max_{Q_1} U_{H,L}(Q_1)$ . Hence, it follows that condition (3.23) for  $x = H$  is equivalent to condition

$$U_{H,H}(Q_1^{H*}) \geq \max\{U_{H,L}(Q_1^{L**}), U_{H,H}(Q^L)\} \quad \text{if } Q_1^{H**} \leq Q^L,$$

---

<sup>4</sup>Recall that  $U_{L,L}(Q_1^{L**}) = \max_{Q_1} U_{L,L}(Q_1)$ .

i.e.

$$Q_1^{H*} \in [Q_1^{H-}, Q_1^{H+}] \quad \text{if } Q_1^{H**} \leq Q^L. \quad (\text{C.5})$$

The proof in more detail:

We first show that no separating equilibrium can exist with  $Q_1^{L*} \neq Q_1^{L**}$ . Assume a separating equilibrium,  $(Q_1^{L*}, Q_2^{L*}, Q_1^{H*}, Q_2^{H*}, \tau^{L*}, \tau^{H*})$ , with  $Q_1^{L*} \neq Q_1^{L**}$ , and  $\mu^* \in \Omega$ , exists. Recall that  $Q_2^{L*} = Q_2^{L**}(Q_1^{L*}, \tau^{L*})$  and  $\tau^{L*} = \tau^{L**}(Q_1^{L*})$ . So the low type's ex ante utility at period 1 in this separating equilibrium is  $U_{L,L}(Q_1^{L*})$ . However, we now show that the strategies cannot constitute an equilibrium as it is always profitable for the low type to deviate from  $Q_1^{L*}$  to  $Q_1^{L**}$  in period 1. Two possible cases can occur,  $Q_1^{H*} \neq Q_1^{L**}$  and  $Q_1^{H*} = Q_1^{L**}$ .

1. Assume  $Q_1^{H*} \neq Q_1^{L**}$ . In this case given the belief system,  $\mu^*$ , the following two equations are true,  $Q_2^{L*}(Q_1^{L**}, \tau^*(Q_1^{L**})) = Q_2^{L**}$  and  $\tau^*(Q_1^{L**}) = \tau^{L**}$ . So the low type's ex ante utility at period 1 given the deviation to  $Q_1^{L**}$  is  $U_{L,L}(Q_1^{L**})$ . From the symmetric information case we know that

$$U_{L,L}(Q_1^{L*}) \leq U_{L,L}(Q_1^{L**}),$$

which however violates the equilibrium condition (3.23).

2. Assume  $Q_1^{H*} = Q_1^{L**}$ . In this case given the belief system,  $\mu^*$ , the following two equations are true,  $Q_2^{L*}(Q_1^{L**}, \tau^*(Q_1^{L**})) = Q_2^{L**}(Q_1^{L**}, \tau^*(Q_1^{L**}))$  and  $\tau^*(Q_1^{L**}) = \tau^{H**}(Q_1^{L**}) \leq \tau^{L**}$ . So the low type's ex ante utility at period 1 given the deviation to  $Q_1^{L**}$  is  $U_{L,H}(Q_1^{L**})$ . Note that as  $\tau^{H**}(Q_1^{L**}) \leq \tau^{L**}$  we have that  $U_{L,H}(Q_1^{L**}) \geq U_{L,L}(Q_1^{L**})$ . Together with the symmetric information analysis we have

$$\begin{aligned} U_{L,H}(Q_1^{L**}) &\geq U_{L,L}(Q_1^{L**}) \\ &\geq U_{L,L}(Q_1^{L*}), \end{aligned}$$

which also violates the equilibrium condition (3.23).

Both results together show that the above stated strategies cannot constitute a separating equilibrium.

Let us divide the remaining analysis into two parts. We first deal with the cases,  $Q_1^{H**} > Q^L$ . We show that no other strategy profile can constitute a separating equilibrium than the

strategy profile given in proposition 3.2. In a second step we show that the stated “symmetric information” strategies indeed constitute a separating equilibrium.

1. From the above analysis we know that the low type’s strategy in period 1 in any separating equilibrium satisfies  $Q_1^{L*} = Q_1^{L**}$ . Assume another separating equilibrium given  $\mu^* \in \Omega$  exists such that  $Q_1^{H*} \neq Q_1^{H**}$ . The high type’s ex ante utility at period 1 in this equilibrium is  $U_{H,H}(Q_1^{H*})$ . We argue that this cannot be an equilibrium as the high type always has an incentive to deviate from  $Q_1^{H*}$  to  $Q_1^{H**}$  in period 1. Given the belief system  $\mu^*$  the regulator’s tax response satisfies  $\tau^*(Q_1^{H**}) = \tau^{H**}$  as  $Q_1^{H**} > Q^L$ . From the symmetric information case we know that

$$U_{H,H}(Q_1^{H*}) \leq U_{H,H}(Q_1^{H**}),$$

which however violates the equilibrium condition (equation (3.23)).

2. We now show that the strategy profile,  $(Q_1^{L**}, Q_2^{L**}, Q_1^{H**}, Q_2^{H**}, \tau^{L**}, \tau^{H**})$  and  $\mu^* \in \Omega$ , indeed constitutes a separating equilibrium.

The low type’s ex ante utility at period 1 in this equilibrium is  $U_{L,L}(Q_1^{L**})$ . The low type does not have an incentive to deviate to any other strategy in period 1. Indeed, the low type deviating from  $Q_1^{L**}$  to any strategy  $Q_1 \leq Q^L$  faces given the belief system  $\mu^*$  the L-Tax-Response,  $\tau^*(Q_1) = \tau^{L**}(Q_1)$ , and supplies,  $Q_2^{L*}(Q_1, \tau^*(Q_1)) = Q_2^{L**}(Q_1, \tau^*(Q_1))$ , in period 2, i.e. his ex ante utility at period 1 is hence  $U_{L,L}(Q_1)$ . In the symmetric information case we have seen that under these circumstances the following inequality holds for all  $Q_1 \leq Q^L$ ,

$$U_{L,L}(Q_1) \leq U_{L,L}(Q_1^{L**}).$$

This means that the low type does not have any incentives to deviate from his stated strategies.

The high type’s ex ante utility at period 1 in this equilibrium is  $U_{H,H}(Q_1^{H**})$ . We now show that also the high type does not have an incentive to deviate to any other strategy in period 1. The high type deviating from  $Q_1^{H**}$  to any strategy  $Q^L < Q_1 \leq Q^H$  faces given  $\mu^*$  the H-Tax-Response,  $\tau^*(Q_1) = \tau^{H**}(Q_1)$ , and supplies,  $Q_2^{H*}(Q_1, \tau^*(Q_1)) = Q_2^{H**}(Q_1, \tau^*(Q_1))$ , in period 2, i.e. his ex ante utility at period 1 is hence  $U_{H,H}(Q_1)$ . In the symmetric information case we have seen that under these circumstances the

following inequality holds for all  $Q^L < Q_1 \leq Q^H$ ,

$$U_{H,H}(Q_1) \leq U_{H,H}(Q_1^{H**}).$$

This means that the high type does not have any incentives to deviate from his stated strategy to any strategy  $Q^L < Q_1 \leq Q^H$ .

We now show that he also does not have an incentive to deviate to any  $Q_1 \leq Q^L$ . The high type deviating from  $Q_1^{H**}$  to any strategy  $Q_1 \leq Q^L$  faces, given  $\mu^*$ , the L-Tax-Response,  $\tau^*(Q_1) = \tau^{L**}(Q_1) \geq \tau^{H**}(Q_1)$ , and supplies,  $Q_2^{H*}(Q_1, \tau^*(Q_1)) = Q_2^{H**}(Q_1, \tau^*(Q_1))$ , in period 2, i.e. his ex ante utility is  $U_{H,L}(Q_1)$ . Together with the symmetric information case we find that under these circumstances the following inequality holds for all  $Q_1 \leq Q^L$ ,

$$\begin{aligned} U_{H,H}(Q_1^{H**}) &\geq U_{H,H}(Q_1) \\ &\geq U_{H,L}(Q_1). \end{aligned}$$

This means that the high type also does not have any incentives to deviate from his stated strategies to any strategy  $Q_1 \leq Q^L$ .

So, the stated strategies indeed constitute a separating equilibrium.

Let us now deal with the cases,  $Q_1^{H**} \leq Q^L$ . We first show that there cannot exist an equilibrium with  $Q_1^{H*} \notin [Q_1^{H-}, Q_1^{H+}]$ . In a second step we find conditions such that the stated strategies constitute an equilibrium given  $\mu^*$ .

Assume a separating equilibrium with  $Q_1^{H*} \notin [Q_1^{H-}, Q_1^{H+}]$  exists. We argue that the high type always has an incentive to deviate from  $Q_1^{H*}$  either to  $Q^L$  or to  $Q_1^{L**}$ . The high type's ex ante utility at period 1 given the equilibrium strategies is  $U_{H,H}(Q_1^{H*})$ . Given the definitions of the thresholds  $Q_1^{H-}$  and  $Q_1^{H+}$  it follows that  $U_{H,H}(Q_1^{H*}) \leq \max\{U_{H,L}(Q_1^{L**}), U_{H,H}(Q^L)\}$ . From this it follows immediately that the high type indeed has an incentive to deviate from  $Q_1^{H*}$  either to  $Q^L$  or to  $Q_1^{L**}$ .

We argue that if  $Q_1^{H*} \in [Q_1^{H-}, Q_1^{H+}]$  the high type does not have an incentive to deviate in period 1. As the high type's utility function  $U_{H,H}(Q_1)$  is concave in  $Q_1$  and is maximized at  $Q_1^{H**} \leq Q^L$  the high type does not have an incentive to deviate from his stated equilibrium strategy to a strategy  $Q_1 > Q^L$  as he subsequently faces the H-Tax-Response

and as  $U_{H,H}(Q_1^{H*}) \geq U_{H,H}(Q^L)$ . Now, recall that given any high type's deviation to a  $Q_1 \leq Q^L$  he subsequently faces the low type's tax and his utility is  $U_{H,L}(Q_1)$ . Now as  $U_{H,H}(Q_1^{H*}) \geq \max U_{H,L}(Q_1)$  the high type indeed does not have an incentive to deviate to any  $Q_1 \leq Q^L$ . Both results together state that the high type does not have an incentive to deviate from his stated equilibrium strategy of period 1.

We now show that the low type does not have an incentive to deviate in period 1 if  $U_{L,L}(Q_1^{L**}) \geq U_{L,H}(Q_1^{H*})$  and he prefers to deviate if this is not the case. The low type deviating to any strategy  $Q_1 \leq Q^L$  with  $Q_1 \neq Q_1^{H*}$  faces his own tax and his utility is  $U_{L,L}(Q_1)$ . As we have seen in the symmetric information case the low type's utility satisfies  $U_{L,L}(Q_1^{L**}) \geq U_{L,L}(Q_1)$ . So this implies that he does not have an incentive to deviate to any  $Q_1 \leq Q^L$  with  $Q_1 \neq Q_1^{H*}$  in period 1. Now note that the stated strategies of proposition 3.2 only constitute an equilibrium if and only if the low type does not have an incentive to deviate to  $Q_1^{H*}$ . This is the same as assuming  $U_{L,L}(Q_1^{L**}) \geq U_{L,H}(Q_1^{H*})$  as he subsequently faces the H-Tax-Response if he deviates to  $Q_1^{H*}$ . ■

### Pooling Equilibria:

**Proof of Lemma 3.5 and Lemma 3.6:** The regulator's problem between period 1 and period 2 is

$$\begin{aligned} \tau^*(Q_1) = \operatorname{argmax}_\tau & [\mu(Q^L | Q_1) U_S^{t=2}(Q_1, Q_2^{L*}(Q_1, \tau), \tau) \\ & + (1 - \mu(Q^L | Q_1)) U_S^{t=2}(Q_1, Q_2^{H*}(Q_1, \tau), \tau)]. \end{aligned}$$

Solving the regulator's problem whenever the belief satisfies either  $\mu(Q^L | Q_1) = 0$  or  $\mu(Q^L | Q_1) = 1$  reduces to the symmetric information tax response problem and gives us straightforwardly the symmetric information tax strategies. So, the only remaining unsolved cases are those with  $\mu(Q^L | Q_1) = \mu$ . In these cases the regulator's problem after observing the supply,  $Q_1$ , is simply

$$\tau^*(Q_1) = \operatorname{argmax}_\tau [\mu U_S^{t=2}(Q_1, Q_2^{L**}(Q_1, \tau), \tau) + (1 - \mu) U_S^{t=2}(Q_1, Q_2^{H**}(Q_1, \tau), \tau)]. \quad (\text{C.6})$$

Recall from the proof of Proposition 3.1 that whenever  $Q_1 \leq \frac{Q^L(r+d(1+r))-r}{r} \leq \frac{Q^H(r+d(1+r))-r}{r}$  the regulator's symmetric information tax responses are  $\tau^{x**}(Q_1) = \tau_{op}(Q_1)$ , for all  $x \in \{L, H\}$ . And the fossil-fuel owner's symmetric information

responses in period 2 are  $Q_2^{x^{**}}(Q_1, \tau_{op}(Q_1)) = \frac{1-\tau_{op}(Q_1)}{2}$ . This in other words means that if  $Q_1 \leq \frac{Q^L(r+d(1+r))-r}{r}$  the regulator does not want the reserves to be exhausted independent of the fossil-fuel owner's type. She even wants both types to supply the same amount in period 2. As she does not distinguish between both types when she observes a  $Q_1 \leq \frac{Q^L(r+d(1+r))-r}{r}$  her best response in the asymmetric information case is<sup>5</sup>

$$\tau^*(Q_1) = \tau^{x^{**}}(Q_1) = \tau_{op}(Q_1), \quad x \in \{L, H\}$$

even if  $0 < \mu < 1$ .

If instead  $Q_1 \geq \frac{Q^L(r+d(1+r))-r}{r}$  the regulator, in the symmetric information case, prefers to set different taxes when facing the different types. Note that for any tax  $\tau \leq \tau^{L^{**}}(Q_1)$  both fossil-fuel owners supply different amounts in period 2. So, if  $0 < \mu < 1$  it is not straightforwardly clear which tax solves equation (C.6).

Before we solve the problem let us first assume that the low type exhaust his reserves and the high type does not exhaust his. If these responses hold for any tax level the regulator's problem (equation (C.6)) becomes

$$\tau_r(Q_1) \equiv \operatorname{argmax}_{\tau} \left\{ \mu \cdot U_S^{t=2}(Q_1, \frac{1-\tau}{2}, \tau) + (1-\mu) \cdot U_S^{t=2}(Q_1, Q^L - Q_1, \tau) \right\}.$$

Solving this problem by the first order condition yields

$$\tau_r(Q_1) = \tau_{op}(Q_1) + \frac{\mu}{(1-\mu)} \frac{(1-2\lambda)(Q^L - Q_1)}{3r - 5\lambda r + d(1+r)} \geq \tau_{op}(Q_1).$$

With the same boundary arguments as in the proof of proposition 3.1 we can now conclude the following statements if  $Q_1 \geq \frac{Q^L(r+d(1+r))-r}{r}$

1. If  $\tau_r(Q_1) \geq \tau_c(Q_1, Q^L)$  the solution to the regulator's problem (equation (C.6)) is the L-Tax-Response

$$\tau^*(Q_1) = \tau_c(Q_1, Q^L).$$

Note that in this situation the following inequalities are true,

---

<sup>5</sup>This in a way results in solving her problem neglecting the probabilities

$$\begin{aligned} \tau^*(Q_1) &= \operatorname{argmax}_{\tau} [\mu U_S^{t=2}(Q_1, Q_2^{L^{**}}(Q_1, \tau), \tau) + (1-\mu) U_S^{t=2}(Q_1, Q_2^{H^{**}}(Q_1, \tau), \tau)] \\ &= \operatorname{argmax}_{\tau} U_S^{t=2}(Q_1, \frac{1-\tau}{2}, \tau). \end{aligned}$$

$\tau^{H**}(Q_1) \leq \tau^* = \tau_c(Q_1, Q^L) = \tau^{L**}(Q_1)$ . The low type exhaust his reserves and the high type does not exhaust his reserves.

2. If  $\tau_c(Q_1, Q^H) \leq \tau_r(Q_1) \leq \tau_c(Q_1, Q^L)$  the solution to the regulator's problem (equation (C.6)) is

$$\tau^*(Q_1) = \tau_r(Q_1).$$

Note that in this situation the following inequalities are true,  $\tau^{H**}(Q_1) \leq \tau^* = \tau_r(Q_1) \leq \tau^{L**}(Q_1)$ . The low type exhaust his reserves and the high type does not exhaust his reserves.

3. If  $\tau_r(Q_1) \leq \tau_c(Q_1, Q^H)$  the solution to the regulator's problem (equation (C.6)) is the high type's symmetric information tax

$$\tau^*(Q_1) = \tau_c(Q_1, Q^H).$$

Note that in this situation the following inequalities are true,

$$\tau^{H**}(Q_1) = \tau^* = \tau_c(Q_1, Q^H) \leq \tau^{L**}(Q_1). \text{ Both types exhaust their reserves.}$$

Also note that the tax,  $\tau_r(Q_1)$ , is an increasing function in  $\mu$ ,  $\frac{\partial}{\partial \mu} \tau_r(Q_1, \mu) \geq 0$ , and takes the following two values at the boundaries,  $\tau_r(Q_1, \mu = 0) = \tau_{op}(Q_1)$  and  $\tau_r(Q_1, \mu = 1) = \infty$ . Hence, the pooling tax,  $\tau^*$ , is also an increasing function in  $\mu$  which satisfies  $\tau^{H**}(Q_1) \leq \tau^* \leq \tau^{L**}(Q_1)$ . All other remaining results (in particular the thresholds) follow directly from these results. ■

**Proof of Proposition 3.3:** Let us first deal with the case  $Q_1^{H**} > Q^L$ . We show that there does not exist a pooling equilibrium. Assume a pooling equilibrium exist, i.e.  $Q_1^* \leq Q^L$ . The high type's ex ante utility at period 1 in this equilibrium is  $U_{H,*}(Q_1^*)$ . We now argue that the high type always has an incentive to deviate to  $Q_1^{H**}$ . Note that given  $\mu^*$  the tax response satisfies  $\tau^*(Q_1^{H**}) = \tau^{H**}$  as  $Q_1^{H**} > Q^L$ . As the P-Tax-Response is higher than the H-Tax-Response,  $\tau^*(Q_1) \geq \tau^{H**}(Q_1)$ , for all  $0 \leq Q_1 \leq Q^H$  and together with the symmetric information analysis we have,

$$U_{H,H}(Q_1^{H**}) \geq U_{H,H}(Q_1^*) \geq U_{H,*}(Q_1^*).$$

However, this violates the equilibrium condition (3.32).



Let us now deal with the cases  $Q_1^{H**} \leq Q^L$ .

We first show that there cannot exist a pooling equilibrium with  $Q_1^{H*} \notin [Q_1^{H-}, Q_1^{H+}]$ . We argue that the high type always has an incentive to deviate from  $Q_1^{H*}$  either to  $Q^L$  or to  $Q_1^{L**}$ . If an equilibrium with  $Q_1^{H*} \notin [Q_1^{H-}, Q_1^{H+}]$  exist the high type's ex ante utility at period 1 in this equilibrium is  $U_{H,*}(Q_1^{H*})$ . Note that  $U_{H,*}(Q_1^{H*}) \leq U_{H,H}(Q_1^{H*})$ . Given the definitions of the thresholds  $Q_1^{H-}$  and  $Q_1^{H+}$  it follows that  $U_{H,H}(Q_1^{H*}) \leq \max\{U_{H,L}(Q_1^{L**}), U_{H,H}(Q^L)\}$ , hence  $U_{H,*}(Q_1^{H*}) \leq \max\{U_{H,L}(Q_1^{L**}), U_{H,H}(Q^L)\}$ . So, the high type does have an incentive to deviate from  $Q_1^{H*}$  either to  $Q^L$  or to  $Q_1^{L**}$ .

Furthermore it is clear that a necessary condition for a strategy profile  $(Q_1^*, Q_2^{L*}, Q_2^{H*}, \tau^*)$  and  $\mu^*$  to constitute a pooling equilibrium is  $U_{H,*}(Q_1^{H*}) \leq \max\{U_{H,L}(Q_1^{L**}), U_{H,H}(Q^L)\}$  as otherwise the high type given  $\mu^*$  has an incentive to deviate from  $Q_1^{H*}$  either to  $Q_1^{L**}$  or to  $Q^L$ .

Note also that a necessary condition for a strategy profile  $(Q_1^*, Q_2^{L*}, Q_2^{H*}, \tau^*)$  to constitute a pooling equilibrium given  $\mu^*$  is  $U_{L,*}(Q_1^*) \geq U_{L,L}(Q_1^{L**})$  as otherwise the low type given  $\mu^*$  has an incentive to deviate to  $Q_1^{L**}$ .

Furthermore, note that as  $U_{L,*}(Q_1) \geq U_{L,L}(Q_1)$  (this is true as the P-Tax-Response is smaller than the L-Tax-Response.) a pooling equilibrium with  $Q_1^* = Q_1^H$  can only exist if and only if a separating equilibrium with the same  $Q_1^{H*} = Q_1^H$  does not exist.

We now argue that the stated strategies constitute a pooling equilibrium given the stated conditions.

1. As the high type's utility function  $U_{H,H}(Q_1)$  is concave in  $Q_1$  and maximized at  $Q_1^{H**} \leq Q^L$  the high type does not have an incentive to deviate from his stated equilibrium strategy to a strategy  $Q_1 > Q^L$  as  $U_{H,*}(Q_1^*) \geq U_{H,H}(Q^L) \geq U_{H,H}(Q_1)$  for all  $Q_1 > Q^L$  (note that he faces his own tax in these cases). Now, recall that given any high type's deviation to a  $Q_1 \leq Q^L$  he subsequently faces the low type's tax and his utility is  $U_{H,L}(Q_1)$ . Now as  $U_{H,*}(Q_1^*) \geq \max U_{H,L}(Q_1)$  the high type also does not have an incentive to deviate to any  $Q_1 \leq Q^L$ . So, he indeed does not have an incentive to deviate from the stated equilibrium strategy  $Q_1^*$  at all.
2. The low type deviating to any strategy  $Q_1 \leq Q^L$  faces his own tax and his utility is  $U_{L,L}(Q_1)$ . As we have seen in the symmetric information case the low type's utilities satisfy  $U_{L,L}(Q_1^{L**}) \geq U_{L,L}(Q_1)$ . So, it follows directly that if  $U_{L,L}(Q_1^{L**}) \leq U_{L,*}(Q_1^*)$  the low type does not have an incentive to deviate from  $Q_1^*$  to any other strategy  $Q_1 \leq Q^L$

in period 1.

■

**Proof of Lemma 3.7:** Assume a pooling equilibrium,  $(Q_1^*, Q_2^{L*}, Q_2^{H*}, \tau^*)$ , and  $\mu^* \in \Omega$ , with  $\mu \in [0, 1]$  exist, i.e. the function  $U_{x,*}(Q_1^*)$  satisfies condition (3.34) for all  $x \in \{L, H\}$ .

Now, as  $\tau^*$  increases in  $\mu$  the function,  $U_{x,*}(Q_1^*)$ , decreases in  $\mu$  for all  $x \in \{L, H\}$ . However, the functions,  $U_{L,L}(Q_1^{L**})$ ,  $U_{H,L}(Q_1^{L**})$ , and  $U_{H,H}(Q_1^L)$  are constant in  $\mu$  as they do not depend on the probability,  $\mu$ , at all. Hence the equilibrium condition (3.34) is harder to met when  $\mu$  increases. The lemma now follows immediately. ■

### Equilibrium Selection:

#### Proof of Lemma 3.8:

Let us first deal with the set of pooling equilibria.

We first show that the stated pooling equilibrium indeed exist whenever  $\mu \leq \mu_T(Q_1^{H**})$ . The fossil-fuel owner  $x \in \{L, H\}$ 's ex ante utility at period 1 is  $U_{x,H}(Q_1^{H**})$  as  $\tau^* = \tau^*(Q_1^{H**}) = \tau^{H**}$ . As additionally  $U_{x,H}(Q_1^{H**}) = \max_{Q_1} U_{x,H}(Q_1)$  both types do not have an incentive to deviate in period 1 from the equilibrium supply to any other supply. And so the stated strategies indeed constitute a pooling equilibrium. We now show that this equilibrium survives the Cho-Kreps' Intuitive Criterion. Assume the regulator always believes that she is facing the high type whatever off equilibrium supply she observes in period 1. In terms of the fossil-fuel owner's interest this is the best belief he can hope for as it generates the lowest possible tax response,  $\tau^{H**}(Q_1)$ . As, even given this belief the fossil-fuel owner does not have an incentive to deviate in period 1 ( $U_{x,H}(Q_1^{H**}) = \max_{Q_1} U_{x,H}(Q_1)$ ) the stated equilibrium survives the Cho-Kreps' Intuitive Criterion.

All other pooling equilibria do not survive the Cho-Kreps' Intuitive Criterion independent of  $\mu$ .

Assume a strategy profile,  $(Q_1^*, Q_2^{L*}, Q_2^{H*}, \tau^*)$ , together with  $\mu^*$  constitute a pooling equilibrium. Note given any pooling equilibrium the following inequality is always true  $U_{x,H}(Q_1^*) \leq U_{x,H}(Q_1^{H**})$  for all  $x \in \{L, H\}$ . We now show that there exist a  $0 \leq Q_1 \leq Q^H$  such that the regulator's reasonable belief is  $\mu^*(Q^L | Q_1) = 0$ . To understand this note the following reasoning. The utility functions,  $U_{H,H}(Q_1)$  and  $U_{L,H}(Q_1)$ , are parabolas which are open downwards and which are compressed at the same rate as  $\frac{\partial^2}{\partial Q_1^2} U_{H,H}(Q_1) = \frac{-2r}{1+r} = \frac{\partial^2}{\partial Q_1^2} U_{L,H}(Q_1)$ . Now calculating the difference in the fossil-fuel owner's ex ante utilities facing the H-Tax-Response and

facing the P-Tax-Response we find that it satisfies  $U_{x,H}(Q_1) - U_{x,*}(Q_1) = (Q^x - Q_1)[\tau^*(Q_1) - \tau^{H**}(Q_1)] \geq 0$  for all  $0 \leq Q_1 \leq Q^H$  and all  $x \in \{L, H\}$ . This implies that this difference is smaller for the low type than for the high type,  $U_{L,H}(Q_1) - U_{L,*}(Q_1) \leq U_{H,H}(Q_1) - U_{H,*}(Q_1)$ , for all  $0 \leq Q_1 \leq Q^H$ . This especially means that  $U_{L,H}(Q_1^*) - U_{L,*}(Q_1^*) \leq U_{H,H}(Q_1^*) - U_{H,*}(Q_1^*)$ . Also recall that the low type's utility is smaller than the high type's utility when facing the same tax response,  $U_{L,H}(Q_1) \leq U_{H,H}(Q_1)$  for all  $Q_1$ . At last let us define two thresholds,  $Q_1^{x,*,-} \equiv \min\{U_{x,H}^{-1}(U_{x,*}(Q_1^*))\}$  and  $Q_1^{x,*,+} \equiv \max\{U_{x,H}^{-1}(U_{x,*}(Q_1^*))\}$  for all  $x \in \{L, H\}$ .<sup>6</sup> Taking the last three results together we can conclude that  $[Q_1^{L,*,-}, Q_1^{L,*,+}] \subsetneq [Q_1^{H,*,-}, Q_1^{H,*,+}]$ . As the intervals are compact sets it follows that there exist an off equilibrium path supply with  $Q_1 \in [Q_1^{H,*,-}, Q_1^{H,*,+}]$  such that  $Q_1 \notin [Q_1^{L,*,-}, Q_1^{L,*,+}]$ . It now follows directly that the low type is worse off deviating from  $Q_1^*$  to this specific  $Q_1$  even if he would face the high type's tax. This implies that the regulator's reasonable belief after observing this off equilibrium path supply is  $\mu^*(Q^L | Q_1) = 0$ . Given the corresponding tax response,  $\tau^{H**}(Q_1)$ , the high type is better off deviating to this off equilibrium path strategy as  $U_{H,H}(Q_1^*) < U_{H,H}(Q_1)$ , and hence the pooling equilibrium does not survive the Cho-Kreps' Intuitive Criterion.

We now deal with the set of separating equilibria. Recall that it is independent of the probability  $\mu$ . Assume a strategy profile,  $(Q_1^{L**}, Q_2^{L**}, Q_1^{H*}, Q_2^{H*}, \tau^{L**}, \tau^{H**})$ , together with  $\mu^*$  constitute a separating equilibrium. One of the necessary conditions for the strategy profile to constitute a separating equilibrium is  $U_{L,L}(Q_1^{L**}) \geq U_{L,H}(Q_1^{H*})$ , because otherwise the low type would have an incentive to mimic the high type. So, as  $U_{L,L}(Q_1^{L**}) \leq U_{L,H}(Q_1^{H**})$  the high type's equilibrium supply satisfies  $Q_1^{H*} \neq Q_1^{H**}$ . We now show that the Cho-Kreps' Intuitive Criterion eliminates all separating equilibria. Let us divide the analysis into two parts. Let us have a look at the cases in which  $U_{L,L}(Q_1^{L**}) > U_{L,H}(Q_1^{H*})$ . As we deal with continuous functions and as  $Q_1^{H**} = \operatorname{argmax}_{Q_1} U_{H,H}(Q_1)$  there exist a  $Q_1$  close to  $Q_1^{H*}$  such that  $U_{L,L}(Q_1^{L**}) > U_{L,H}(Q_1)$  and  $U_{H,H}(Q_1) > U_{H,H}(Q_1^{H*})$ . Hence, the regulator's reasonable belief given this specific  $Q_1$  is  $\mu^*(Q^L | Q_1) = 0$  and given the corresponding tax response the high type has an incentive to deviate from  $Q_1^{H*}$  to this particular  $Q_1$ . To reason in the same manner as we did for the pooling equilibrium, note the following. Let us first define two thresholds  $Q_1^{x,ic-} \equiv \min U_{x,H}^{-1}(U_{x,x}(Q_1^{x*}))$  and  $Q_1^{x,ic+} \equiv \max U_{x,H}^{-1}(U_{x,x}(Q_1^{x*}))$  for all  $x \in \{L, H\}$ . The corresponding sets are not empty as  $Q_1^{H**} \in [Q_1^{x,ic-}, Q_1^{x,ic+}]$  for all  $x \in \{L, H\}$ . Furthermore it is also true that  $Q_1^{H*} \in \{Q_1^{H,ic-}, Q_1^{H,ic+}\}$ . We now show that the low type's

---

<sup>6</sup>Note that the interval,  $[Q_1^{x,*-}, Q_1^{x,*+}]$  is not empty as  $Q_1^{H**} \in [Q_1^{x,*-}, Q_1^{x,*+}]$  for all  $x \in \{L, H\}$ .

interval is a strict subset of the high type's interval,  $[Q_1^{L,ic-}, Q_1^{L,ic+}] \subsetneq [Q_1^{H,ic-}, Q_1^{H,ic+}]$ . To understand this note that, i) the following inequality is true,  $U_{L,L}(Q_1^{L**}) > U_{L,H}(Q_1^{H*})$ , and ii) the high type's equilibrium supply,  $Q_1^{H*}$ , naturally satisfies  $Q_1^{H*} \in U_{L,H}^{-1}(U_{L,H}(Q_1^{H*}))$ . These facts imply that  $Q_1^{H*} \notin [Q_1^{L,ic-}, Q_1^{L,ic+}]$ . However, it is also true that  $Q_1^{H*} \in \{Q_1^{H,ic-}, Q_1^{H,ic+}\}$ . Now, as the low and the high type's intervals are compact sets there exist an off equilibrium path supply  $Q_1 \in [Q_1^{H,ic-}, Q_1^{H,ic+}] \setminus [Q_1^{L,ic-}, Q_1^{L,ic+}]$ , and hence the separating equilibrium does not survive the Cho-Kreps' Intuitive Criterion. Let us now deal with the cases in which  $U_{L,L}(Q_1^{L**}) = U_{L,H}(Q_1^{H*})$ . Note that as  $U_{H,H}(Q_1)$  is a parabola there exist a supply  $Q_1 \neq Q_1^{H**}$  such that  $U_{H,H}(Q_1) = U_{H,H}(Q_1^{H*})$ . Applying the above techniques to this  $Q_1$  again implies that the separating equilibrium does not survive the Cho-Kreps' Intuitive Criterion. ■

### C.3 Welfare Analysis

**Proof of Proposition 3.4:** The  $x \in \{L, H\}$  type's ex ante utility at period 1 given this pooling equilibrium is  $U_{x,H}(Q_1^{H**})$ . Recall that  $U_{x,H}(Q_1^{H**}) = \max_{Q_1} U_{x,H}(Q_1)$ . As the H-Tax-Response is smaller than the P-Tax-Response which is again smaller than the L-Tax-Response,  $\tau^{H**}(Q_1) \leq \tau^*(Q_1) \leq \tau^{L**}(Q_1)$ , the following inequalities are also true,  $U_{x,H}(Q_1) \geq U_{x,*}(Q_1) \geq U_{x,L}(Q_1)$  for all  $0 \leq Q_1 \leq Q^H$  and all  $x \in \{L, H\}$ . Hence, the statement of this proposition follows directly. ■

**Proof of Proposition 3.5:** Note that the high type's strategies coincide in both equilibria. So, the comparison of the regulator's ex ante expected welfare at period 0 only depends on the low type's strategies, i.e.

$$\begin{aligned} \mathbb{E}U_{asy}(Q^L, Q^H) &> \mathbb{E}U_{sym}(Q^L, Q^H) \\ \Leftrightarrow U_S(Q_1^{H**}, Q^L - Q_1^{H**}, \tau^{H**}) &> U_S(Q_1^{L**}, Q_2^{L**}, \tau^{L**}). \end{aligned} \tag{C.7}$$

Recall that the regulator's welfare is

$$U_S(Q_1, Q_2, \tau) = \lambda U_F(Q_1, Q_2, \tau) + (1-\lambda)[U_C(Q_1, Q_2, \tau) + \Upsilon(Q_1, Q_2, \tau)] - C(Q_1) - \frac{1}{r}C(Q_1 + Q_2).$$

From proposition 3.4 we know that

$$U_F(Q_1^{H**}, Q^L - Q_1^{H**}, \tau^{H**}) \geq U_F(Q_1^{L**}, Q_2^{L**}, \tau^{L**}).$$

It follows that a sufficient condition for equation (C.7) to be true is,<sup>7</sup>

$$U_{(S-\lambda F)}(Q_1^{H**}, Q^L - Q_1^{H**}, \tau^{H**}) > U_{(S-\lambda F)}(Q_1^{L**}, Q_2^{L**}, \tau^{L**}). \quad (C.8)$$

Our aim is now to show that there exist pairs  $Q^L$  and  $Q^H$ , with  $Q^L < Q^H$ , such that they satisfy equation (C.8).

Let  $Q^L = \frac{1}{2}$ . This implies that  $Q_1^{L**} = \frac{(1+r)-2Q^L}{2r} = Q^L = \frac{1}{2}$ . We now show that there exist  $\frac{1}{2} < Q^H \leq Q_k^2$  which satisfy the above equation. Recall that  $Q_1^{H**} = \frac{(1+r)-2Q^H}{2r} \leq \frac{1}{2}$  if  $\frac{1}{2} \leq Q^H \leq Q_k^2$ .

We analyze the problem in two steps. First assume  $\lambda = 0$ . This implies

$$U_{(S-0F)}(Q_1^{L**}, Q_2^{L**}, \tau^{L**}) = \frac{1}{8}[1 - d - \frac{d}{r}]$$

and

$$\begin{aligned} U_{(S-0F)}(Q_1^{H**}, Q^L - Q_1^{H**}, \tau^{H**}) &= \frac{1}{2}(Q_1^{H**})^2 + \frac{1}{2} \frac{1}{1+r} (\frac{1}{2} - Q_1^{H**})^2 \\ &\quad + \frac{1}{1+r} (1 - 2Q_1^{H**} + 2Q_1^{H**}) (\frac{1}{2} - Q_1^{H**}) - \frac{d}{2}(Q_1^{H**})^2 - \frac{d}{8r} \\ &= \frac{1}{2}(Q_1^{H**})^2 + \frac{1}{2} \frac{1}{1+r} (\frac{1}{2} - Q_1^{H**})^2 \\ &\quad + \frac{1}{2} \frac{1}{1+r} (1 - 2Q_1^{H**})(2Q_1^{H**}(1+r) - r) - \frac{d}{2}(Q_1^{H**})^2 - \frac{d}{8r}. \end{aligned}$$

Now calculating the difference,

$U_{(S-0F)}(Q_1^{H**}, Q^L - Q_1^{H**}, \tau^{H**}) - U_{(S-0F)}(Q_1^{L**}, Q_2^{L**}, \tau^{L**})$ , and rearranging gives us

$$\begin{aligned} &U_{(S-0F)}(Q_1^{H**}, Q^L - Q_1^{H**}, \tau^{H**}) - U_{(S-0F)}(Q_1^{L**}, Q_2^{L**}, \tau^{L**}) = \\ &\quad - 4(2 + 3r + d(1+r))(Q_1^{H**})^2 + 4(1+4r)Q_1^{H**} + d(1+r) - 5r. \end{aligned}$$

Note that  $U_{(S-0F)}(Q_1^{H**}, Q^L - Q_1^{H**}, \tau^{H**}) - U_{(S-0F)}(Q_1^{L**}, Q_2^{L**}, \tau^{L**}) = 0$  if  $Q_1^{H**} = \frac{1}{2}$  as the strategies in each period coincide.

---

<sup>7</sup> $U_{(S-\lambda F)}(Q_1, Q_2, \tau) \equiv U_S(Q_1, Q_2, \tau) - U_F(Q_1, Q_2, \tau)$

We now show that there exist a threshold  $Q^R \leq \frac{1}{2}$  such that all  $Q_1^{H**} \in [Q^R, \frac{1}{2})$  satisfy equation (C.8). This implies that there exist a threshold  $Q^S$  such that all  $Q^H \in (\frac{1}{2}, Q^S]$  satisfy equation (C.8). Note that to show this it is sufficient to show that the derivative of the term,  $U_{(S-0F)}(Q_1^{H**}, Q^L - Q_1^{H**}, \tau^{H**}) - U_{(S-0F)}(Q_1^{L**}, Q_2^{L**}, \tau^{L**})$ , with respect to  $Q_1^{H**}$  is negative at  $Q_1^{H**} = \frac{1}{2}$ . This is indeed sufficient as the function is a parabola which is open downwards. Taking the derivative yields  $-8(2 + 3r + d(1 + r))Q_1^{H**} + 4(1 + 4r)$ . Plugging in  $Q_1^{H**} = \frac{1}{2}$  gives us

$$-4(2 + 3r + d(1 + r)) + 4 + 16r = -8 - 12r - 4d(1 + r) + 4 + 16r = 4r - 4 - 4d(1 + r) \leq 0.$$

The cases  $0 < \lambda \leq \frac{1}{2}$  follow directly from the above analysis. To see this note the following

$$\begin{aligned} & U_{(S-0F)}(Q_1^{H**}, Q^L - Q_1^{H**}, \tau^{H**}) - U_{(S-0F)}(Q_1^{L**}, Q_2^{L**}, \tau^{L**}) \geq 0 \\ \Leftrightarrow & C(Q_1^{L**}) - C(Q_1^{H**}) \geq U_C(Q_1^{L**}, Q_2^{L**}, \tau^{L**}) + \Upsilon(Q_1^{L**}, Q_2^{L**}, \tau^{L**}) \\ & - U_C(Q_1^{H**}, Q^L - Q_1^{H**}, \tau^{H**}) + \Upsilon(Q_1^{H**}, Q^L - Q_1^{H**}, \tau^{H**}). \end{aligned}$$

As the left hand side of the last inequality is positive (this is true as  $Q_1^{L**} \geq Q_1^{H**}$ ) and as  $1 - \lambda \leq 1$  the above inequality also implies

$$\begin{aligned} C(Q_1^{L**}) - C(Q_1^{H**}) & \geq (1 - \lambda)[U_C(Q_1^{L**}, Q_2^{L**}, \tau^{L**}) + \Upsilon(Q_1^{L**}, Q_2^{L**}, \tau^{L**}) \\ & - U_C(Q_1^{H**}, Q^L - Q_1^{H**}, \tau^{H**}) + \Upsilon(Q_1^{H**}, Q^L - Q_1^{H**}, \tau^{H**})]. \end{aligned}$$

This in turn is the same as

$$U_{(S-\lambda F)}(Q_1^{H**}, Q^L - Q_1^{H**}, \tau^{H**}) - U_{(S-\lambda F)}(Q_1^{L**}, Q_2^{L**}, \tau^{L**}) \geq 0.$$

■







# Bibliography

- Akerberg, D. & Botticini, M. (2002), 'Endogenous matching and the empirical determinants of contract form', *Journal of Political Economy* **110**, 564–591.
- Aggarwal, R. & Samwick, A. (2002), 'The other side of the trade-off: The impact of risk on executive compensation - a reply'.
- Araujo et al. (2007), 'The trade-off between incentives and endogenous risk', *Revista de Econometria* **27**, 193–198.
- Armstrong, M. (1996), 'Multiproduct nonlinear pricing', *Econometrica* **64**, 51–75.
- Austen-Smith, D. & Feddersen, T. J. (2006), 'Deliberation, preference uncertainty and voting rules', *The American Political Science Review* **100**, 209–217.
- Bauer, N., Hilaire, J. & Bertram, C. (2014), 'The calm before the storm - what happens to CO<sub>2</sub> emissions before their price starts to increase?', *Potsdam Institute for Climate Impact Research* .
- Bentley, R. W. (2002), 'Global oil and gas depletion: an overview', *Energy Policy* **30**, 189–205.
- Berle, A. & Means, G. (1932), 'The modern corporation and private property', *New York: Macmillan* .
- Biais, B., Martimort, D. & Rochet, J.-C. (2000), 'Competing mechanisms in a common value environment', *Econometrica* **68**, 799–838.
- Budde, J. & Kräkel, M. (2011), 'Limited liability and the risk-incentive relationship', *Journal of Economics* **102**, 97–110.
- Bushman, R. M. et al. (1996), 'CEO compensation: The role of individual performance evaluation', *Journal of Accounting and Economics* **21**, 161–193.

- Chakravorty, U., Magné, B. & Moreaux, M. (2006), ‘A hotelling model with a ceiling on the stock of pollution’, *Journal of Economic Dynamics and Control* **30**, 2875–2904.
- Chakravorty, U., Moreaux, M. & Tidball, M. (2008), ‘Ordering the extraction of polluting nonrenewable resources’, *American Economic Review* **98**, 1128–44.
- Core, J. & Guay, W. (1999), ‘The use of equity grants to manage optimal equity incentive levels’, *Journal of Accounting and Economics* **28**, 151–184.
- Coughlan, P. J. (2000), ‘In defense of unanimous jury verdicts: mistrials, communication and strategic voting’, *The American Political Science Review* **94**, 375–393.
- Dahl, C. & Yucel, M. (1991), ‘Testing alternative hypotheses of oil producer behavior’, *The Energy Journal* **12**, 117–138.
- Deimen, I., Ketelaar, F. & Le Qument, M. T. (2014), ‘Consistency and communication in committees’, *University of Bonn, mimeo* .
- Demski, J. & Dye, R. (1999), ‘Risk, return, and moral hazard’, *Journal of Accounting Research* **37**, 27–55.
- Doraszelski, U., Gerardi, D. & Squintani, F. (2003), ‘Communication and voting with doubled sided information’, *Contributions to Theoretical Economics* **3**, 1–41.
- Duggan, J. & Martinelli, C. (2001), ‘A bayesian model of voting in juries’, *Games and Economic Behavior* **37**, 259–294.
- Esary, J. D. & Proschan, F. (1972), ‘Relationship among some concepts of bivariate dependence’, *The Annals of Mathematical Statistics* **43**, 651–655.
- Esary, J. D., Proschan, F. & Walkup, D. W. (1967), ‘Association of random variables with applications’, *The Annals of Mathematical Statistics* **38**, 1466–1474.
- Ezzati, A. (1976), ‘Future OPEC price and production strategies as affected by its capacity to absorb oil revenues’, *European Economic Review* **8**, 107–138.
- Feddersen, T. J. & Pesendorfer, W. (1998), ‘Convicting the innocent: the inferiority of unanimous jury verdicts under strategic voting’, *The American Political Science Review* **92**, 23–35.

- Garcia, D. (2014), ‘Optimal contracts with privately informed agents and active principals’, *Journal of Corporate Finance* **29**, 695–709.
- Garen, J. (1994), ‘Executive compensation and principal-agent theory’, *The Journal of Political Economy* **102**, 1175–1199.
- Gaudet, G., Lassere, P. & Long, N. V. (1995), ‘Optimal resource royalties with unknown and temporally independent extraction cost structures’, *International Economic Review* **36**, 715–49.
- Gerardi, D. (2000), ‘Jury verdicts and preference diversity’, *The American Political Science Review* **94**, 395–406.
- Gerardi, D. & Yariy, L. (2007), ‘Deliberative voting’, *Journal of Economic Theory* **134**, 317–338.
- Greene, W. H. (1993), ‘Econometric analysis’, *Prentice Hall International* .
- Griffin, J. M. (1985), ‘OPEC behavior: A test of alternative hypotheses’, *American Economic Review* **75**, 954–63.
- Hellwig, M. (2009), ‘A reconsideration of the Jensen-Meckling model of outside finance’, *Journal of Financial Intermediation* **18**, 495–525.
- Hellwig, M. & Schmidt, K. (2002), ‘Discrete-time approximations of the Holmström-Milgrom Brownian-motion model of intertemporal incentive provision’, *Econometrica* **70**, 2225–2264.
- Holmström, B. & Milgrom, P. (1987), ‘Aggregation and linearity in the provision of intertemporal incentives’, *Econometrica* **55**, 303–328.
- Holmström, B. & Milgrom, P. (1994), ‘The firm as an incentive system’, *American Economic Association* **84**, 972–991.
- Hotelling, H. (1931), ‘The economics of exhaustible resources’, *Journal of Political Economy* **39**, 137.
- Hummel, P. (2012), ‘Deliberation in large juries with diverse preferences’, *Public Choice* **150**, 595–608.

- Hung, N. M., Poudou, J.-C. & Thomas, L. (2006), ‘Optimal resource extraction contract with adverse selection’, *Resources Policy* **31**, 78–85.
- IEA (2010), ‘World energy outlook 2010’, *International Energy Agency (IEA)* .
- IEA (2011), ‘CO<sub>2</sub> emissions from fuel combustion’, *International Energy Agency (IEA)* .
- IEA (2012), ‘World energy outlook 2012’, *International Energy Agency (IEA)* .
- Inderst, R. & Müller, H. (2010), ‘CEO replacement under private information’, *Review of Financial Studies* **23**, 2935–2969.
- Ing, J. (2012), ‘The impact of commitment on nonrenewable resources management with asymmetric information on costs’, *HAL, working paper* .
- IPCC (2013), ‘Climate change 2013: The physical science basis’, *Cambridge University Press* .
- Ittner et al. (1997), ‘The choice of performance measures in annual bonus contracts’, *The Accounting Review* **72**, 231–255.
- Jebjerg, L. & Lando, H. (1997), ‘Regulating a polluting firm under asymmetric information’, *Environmental and Resource Economics* **10**, 267–284.
- Jones, C. T. (1990), ‘OPEC behaviour under falling prices: Implications for cartel stability’, *The Energy Journal* **11**, 117–130.
- Lafontaine, F. (1992), ‘Agency theory and franchising: Some empirical results’, *The RAND Journal of Economics* **23**, 263–283.
- Lafontaine, F. & Slade, M. (1998), ‘Incentive contracting and the franchise decision’, *NBER working paper 6544* .
- Lafontaine, F. & Slade, M. (2007), ‘Vertical integration and firm boundaries: The evidence’, *Journal of Economic Literature* **45**, 629–685.
- Laherrere, J. (2013), ‘Shortened world oil & gas production forecasts 1900-2100’, *ASPO France* .

- Lambert, R. & Larcker, D. (1987), ‘An analysis of the use of accounting and market measures of performance in executive compensation contracts’, *Journal of Accounting Research* **25**, 85–125.
- Le Quement, M. (2012), ‘Communication compatible voting rules’, *Theory and Decision* **74**, 479–507.
- Le Quement, M. & Yokeeswaran, V. (2015), ‘Subgroup deliberation and voting’, *Social Choice and Welfare* .
- Loderer, C. (1985), ‘A test of the OPEC cartel hypothesis: 1974-1983’, *The Journal of Finance* **40**, 991–1006.
- Long, N. (2011), ‘Dynamic games in the economics of natural resources: A survey’, *International Review of Economics* **1**, 115–148.
- MacAvoy, P. W. (1982), ‘Crude oil prices as determined by OPEC and market fundamentals’, *Ballinger Publishing Company, Cambridge, MA* .
- Malik, A. S. (1991), ‘Permanent versus interim regulations: A game-theoretic analysis’, *Journal of Environmental Economics and Management* **21**, 127–139.
- Meirowitz, A. (2002), ‘Informative voting and condorcet jury theorems with a continuum of types’, *Social Choice and Welfare* **19**, 219–236.
- Meirowitz, A. (2007), ‘In defense of exclusionary deliberation: Communication and voting with private beliefs and values’, *Journal of Theoretical Politics* **19**, 301–328.
- Osmundsen, P. (1995), ‘Taxation of petroleum companies possessing private information’, *Resource and Energy Economics* **17**, 357–377.
- Osmundsen, P. (1998), ‘Dynamic taxation of non-renewable natural resources under asymmetric information about reserves’, *Canadian Journal of Economics* **31**, 933–951.
- Oyer, P. & Schaefer, S. (2004), ‘Why do some firms give stock options to all employees?: An empirical examination of alternative theories’, *NBER Working Paper 10222* .
- Prendergast, C. (2002), ‘The tenuous trade-off between risk and incentives’, *Journal of Political Economy* **110**, 1071–1102.

- Raith, M. (2003), ‘Competition, risk and managerial incentives’, *American Economic Review* **93**, 1425–1436.
- Rochet, J.-C. & Choné, P. (1998), ‘Ironing, sweeping, and multidimensional screening’, *Econometrica* **66**, 783–826.
- Saure, P. (2010), ‘Overreporting oil reserves’, *Swiss National Bank, working paper* .
- Serfes, K. (2005), ‘Risk sharing vs. incentives: Contract design under two-sided heterogeneity’, *Economics Letters* **88**, 343–349.
- Sung, Y. (2005), ‘Optimal contracts under moral hazard and adverse selection: a continuous time approach’, *Review of Financial Studies* **18**, 1021–1073.
- Tahvonen, O. (1997), ‘Fossil fuels, stock externalities, and backstop technology’, *Canadian Journal of Economics* **30**, 855–74.
- Ulph, A. & Ulph, D. (1994), ‘The optimal time path of a carbon tax’, *Oxford Economic Papers* **46**, 857–68.
- Van der Ploeg, F. & Withagen, C. (2012), ‘Too much coal, too little oil’, *Journal of Public Economics* **96**, 62–77.
- Van Weelden, R. (2008), ‘Deliberation rules and voting’, *Quarterly Journal of Political Science* **3**, 83–88.
- Verleger, Philip K., J. (1982), ‘The determinants of official OPEC crude prices’, *The Review of Economics and Statistics* **64**, 177–82.
- Weinschenk, P. (2014), ‘Sharecropping and performance: A reexamination of the Marshallian hypothesis’, *MPI Bonn, mimeo* .
- Wirl, F. (1994), ‘Pigouvian taxation of energy for flow and stock externalities and strategic, noncompetitive energy pricing’, *Journal of Environmental Economics and Management* **26**, 1–18.
- Wolinsky, A. (2002), ‘Eliciting information from multiple experts’, *Games and Economic Behavior* **41**, 141–160.

Wright, D. (2004), 'The risk and incentive trade-off in the presence of heterogeneous managers', *Journal of Economics* **83**, 209–233.

Yermack, D. L. (1995), 'Do corporations award CEO stock options effectively', *Journal of Financial Economics* **39**, 237–269.