

# Regulating a Monopolist With Uncertain Costs Without Transfers\*

Manuel Amador

Federal Reserve Bank of Minneapolis

University of Minnesota

Kyle Bagwell<sup>†</sup>

Stanford University

April 11, 2019

## Abstract

We analyze the [Baron and Myerson \(1982\)](#) model of regulation under the restriction that transfers are infeasible. To do this, we extend the Lagrangian approach to delegation problems of [Amador and Bagwell \(2013\)](#) to include an ex post participation constraint that allows for the possible exclusion of some types. We report sufficient conditions under which optimal regulation takes the simple and common form of price-cap regulation. We identify families of demand and distribution functions and welfare weights that satisfy our sufficient conditions. We also report conditions under which the optimal price cap is set at a level such that no types are excluded. Using a linear demand example, we show that exclusion of higher cost types can be optimal when these conditions fail to hold. Our analysis also can be used to provide conditions for the optimality of price-cap regulation when an ex post participation constraint is present and exclusion is infeasible.

## 1 Introduction

The optimal regulatory policy for a monopolist is influenced by many considerations, including the possibility of private information, the objective of the regulator, and the feasibility

---

\*The views expressed herein are those of the authors and not necessarily those of the Federal Reserve Bank of Minneapolis or the Federal Reserve System.

<sup>†</sup>Previous draft: November 15, 2016. We would like to thank Mark Armstrong, Yeon-Koo Che, Joe Harrington, Roger Noll, Alessandro Pavan, Michael Riordan, Peter Troyan, Robert Wilson, Frank Wolak and seminar participants at Columbia, Nottingham, Penn, Ryerson, SMU, Vanderbilt and USC for helpful discussions. We thank Guillaume Sublet for providing excellent research assistance. Manuel Amador acknowledges NSF support under award number 0952816. Kyle Bagwell thanks the Center for Advanced Studies in the Behavioral Sciences at Stanford for hospitality and support as a Fellow during 2014-15. Corresponding author, email: amador.manuel@gmail.com; fax: 612-204-5515.

and efficiency of transfers. Simple solutions obtain in some settings. For example, in the textbook case of a single-product monopolist with constant marginal cost and a positive fixed cost, with all costs commonly known, a regulator that maximizes aggregate social surplus obtains the “first-best” (“second-best”) solution by setting price equal to marginal (average) cost when transfers are feasible and efficient (are infeasible). In other settings, however, optimal regulation can take more subtle forms. [Armstrong and Sappington \(2007\)](#) survey the nature of optimal regulation in different settings and discuss as well the design of practical policies, such as price-cap regulation, that are frequently observed in practice. As they emphasize, an important question is whether practical policies perform well in realistic settings where private information may be present and transfer instruments may be limited.

In a seminal paper, [Baron and Myerson \(1982\)](#) consider the optimal regulation of a single-product monopolist with private information about its costs of production. In their model, a regulatory policy indicates, for every possible cost type, whether the monopoly is allowed to produce at all and, if so, the output and corresponding price that it selects and the transfer from consumers that it receives (where a negative transfer is a tax). A regulatory policy is feasible if it is incentive compatible and satisfies an ex post participation constraint. The regulator chooses over feasible regulatory policies to maximize a weighted social welfare function that weighs consumer surplus no less heavily than producer surplus.<sup>1</sup>

In a standard version of the Baron-Myerson model, the monopolist incurs a commonly known and non-negative fixed cost and is privately informed as to the level of its constant marginal cost, where the monopolist’s marginal cost has a continuum of possible types and is drawn from a commonly known distribution function. If the regulator gives greater welfare weight to consumer surplus, then the optimal regulatory policy defines a non-decreasing price schedule for active types with a positive mark up for all but the lowest cost type. By comparison, if the regulator were to maximize aggregate social surplus, then as [Loeb and Magat \(1979\)](#) observe the optimal regulatory policy would achieve a first-best outcome, with price equal to marginal cost for all active types and transfers set so that the monopolist receives all the surplus. In any case, production is permitted only for types such that consumer surplus under the optimal pricing rule weakly exceeds the fixed cost of production.

In this paper, we characterize optimal regulatory policy in the Baron-Myerson model with constant marginal costs when transfers are infeasible. Our no-transfers assumption contrasts sharply with Baron and Myerson’s assumption that all (positive and negative) transfers are available. We motivate our no-transfers assumption in three ways. First, as is commonly

---

<sup>1</sup>An alternative approach is developed by [Laffont and Tirole \(1993, 1986\)](#). They assume that the regulator maximizes aggregate social surplus and that transfers are inefficient (i.e., transfers entail a social cost of funds). Under this approach, consumers incur a cost in excess of one dollar for every dollar that is received as a transfer by the monopolist.

observed, regulators often do not have the authority to explicitly tax or pay subsidies.<sup>2</sup> Second, while transfers from consumers to firms may also be achieved via access fees in two-part tariff schemes, the scope for such transfers may be limited in practice, particularly when universal service is sought and consumers are heterogeneous.<sup>3</sup> Finally, in other settings, the scope for a positive access fee may be limited by the possibility of consumer arbitrage, while the scope for a negative access fee may be limited by the prospect of strategic consumer behavior designed to capture “sign-up” bonuses. In view of these considerations, we remove the traditional assumption that all transfers are available and consider the opposite case in which all transfers are infeasible. Specifically, we assume that the regulated firm is restricted to a uniform price (i.e., linear pricing).<sup>4</sup> As our main finding, we report sufficient conditions under which price-cap regulation emerges as the optimal regulatory policy.

As mentioned above, price-cap regulation is a common form of regulation. The appeal of price-cap regulation is often associated with the incentive that it gives to the regulated firm to invest in endogenous cost reduction.<sup>5</sup> By contrast, we establish conditions for the optimality of price-cap regulation in a model in which costs are private and exogenous. We note further that our no-transfers assumption is critical: price-cap regulation is not optimal in the standard Baron-Myerson model with transfers. Our finding thus indicates that this practical regulatory policy may perform not just well but optimally when a regulator faces a privately informed monopolist and transfers are infeasible.

To develop this finding, we consider a “regulator’s problem” in which the regulator chooses a menu of permissible outputs, with the understanding that the output choice intended for a monopolist with a given cost type must be the best choice for the monopolist relative to all other permitted output choices. In addition to this incentive compatibility constraint, the regulator faces an ex post participation, or individual rationality (IR), constraint: if the regulator seeks a positive output from a monopolist with a given cost type, then the monopolist must earn more by producing this output than by shutting down and avoiding the non-negative fixed cost of production. Importantly, the regulator may choose a

---

<sup>2</sup>For further discussion, see, e. g., [Armstrong and Sappington \(2007, p. 1607\)](#), [Baron \(1989, p. 1351\)](#), [Church and Ware \(2000, p. 840\)](#), [Joskow and Schmalensee \(1986, p. 5\)](#), [Laffont and Tirole \(1993, p. 130\)](#) and [Schmalensee \(1989, p. 418\)](#).

<sup>3</sup>As [Laffont and Tirole \(1993, p. 151\)](#) explain, “optimal linear pricing is a good approximation to optimal two-part pricing when there is concern that a nonnegligible fixed premium would exclude either too many customers or customers with low incomes whose welfare is given substantial weight in the social welfare function.”

<sup>4</sup>In this respect, we follow the lead of [Schmalensee \(1989\)](#), who also examines a regulatory model with linear pricing schemes. [Schmalensee \(1989, p. 418\)](#) provides additional motivation for the practical relevance of linear pricing schemes in regulatory settings.

<sup>5</sup>For further discussion, see, for example, [Armstrong and Sappington \(2007, p. 1608\)](#) and the references cited therein.

menu of permissible outputs such that, for some cost types, the monopolist elects to produce zero output and thus earn a profit of zero. In other words, and as in the original Baron-Myerson model, the regulator may design the regulatory policy so as to “exclude” some cost types from production.

The IR constraint plays an important role in our analysis. If we were to ignore this constraint, then the regulator’s problem would take the form of a traditional delegation problem and fit into the framework of [Amador and Bagwell \(2013\)](#). We could then use the sufficiency theorems developed in that paper and provide conditions under which a simple price cap (i.e., a quantity floor) is optimal. We show, however, that the IR constraint in fact would be violated for higher cost types when this simple price cap is used. Thus, this solution is not feasible for the regulator problem that we analyze here.

We consider instead a price-cap allocation where the cap is placed at a price level such that a threshold cost type earns zero profit and is thus indifferent to shut down. No exclusion occurs if the threshold cost type corresponds to the highest cost type in the full support of possible cost types, while exclusion occurs when the threshold cost type falls below the highest possible cost type. Within the set of non-excluded cost types, higher cost types pool at the price cap, whereas lower cost types may select their monopoly prices. It is also possible that the price cap falls below the monopoly price for the lowest possible cost type, in which case all non-excluded cost types pool at the price cap. The central task of our analysis is to identify sufficient conditions under which the described price cap with possible exclusion is optimal. We also seek to determine sufficient conditions that indicate when actual exclusion does or does not occur.

To establish our results, we proceed in three main steps. First, we consider the “regulator’s truncated problem,” wherein the regulator allocates production for cost types up to an exogenous upper-bound cost type and is not allowed to exclude any types in this truncated set. The upper-bound cost type can be fixed at any value that is above the lowest possible cost type and at or below the highest possible cost type in the full support. We then obtain sufficient conditions under which the optimal allocation for the regulator’s truncated problem is a price cap set at a level such that the upper-bound cost type earns zero profit and is thus indifferent between producing or not. Second, we argue that this truncated allocation remains feasible when extended to the full support of possible costs if cost types above the upper-bound cost type are excluded (assigned zero output). Finally, we characterize the optimal level of exclusion. This exercise amounts to a single variable optimization problem defined over the upper-bound, or threshold, cost type.

Our first proposition establishes a general set of sufficient conditions under which the described price-cap allocation solves the regulator’s truncated problem. The conditions are

defined in terms of general relationships between functions that describe the regulator's welfare, the monopolist's profit and the distribution of cost types, respectively. We then provide a second proposition which establishes that, if the sufficient conditions for our first proposition hold for any upper-bound cost type, then a price-cap allocation with possible exclusion is optimal within the set of all feasible allocations for the regulator's problem. A key ingredient in making this argument is that the optimal price-cap allocation is such that the threshold cost type is indifferent to shut down.

We also provide several results that facilitate the application of our propositions. Three approaches are developed. First, we show that our sufficient conditions hold if the density is non-decreasing over the full support and if a "relative concavity" condition holds that concerns the relative curvature of the consumer surplus and profit functions, with each expressed as a function of quantity. The relative concavity condition is more likely to hold when ratio of the concavity of the consumer surplus function to that of the profit function is higher and when the welfare weight attached to profit is lower. Second, we identify a family of demand functions under which the sufficient conditions for our propositions hold if a simple inequality is satisfied. The inequality condition holds when the density is non-decreasing over the full support, but it can hold as well when the density is decreasing over part or all of the full support. To illustrate the power of this approach, we show that the family includes linear demand, constant elasticity demand and log demand functions, and we derive and interpret the corresponding inequality condition for each of these examples. The third approach is to check the sufficient conditions for our propositions directly. We illustrate this approach for an example with an exponential demand function.

Finally, we turn to the third step of our analysis and identify conditions under which actual exclusion does or does not occur, respectively. Our third proposition establishes that no exclusion is optimal under a general set of conditions; specifically, if the density is non-decreasing over the full support and the consumer surplus function is weakly concave, and if the sufficient conditions for our first proposition hold for any upper-bound cost type, then the optimal regulatory policy entails no exclusion and a price cap set at a price level such that the IR constraint for the highest cost type is binding. Thus, under these conditions, optimal regulation takes the form of a standard second-best price cap that delivers zero profit for the highest cost type. We note that the consumer surplus function is weakly concave in quantity for the log demand and constant elasticity demand examples.

We also analyze the linear demand example. The consumer surplus function associated with this demand function is strictly convex, and so our third proposition cannot be applied. In our fourth and final proposition, we show that, if the distribution of cost types is uniform, the social planner maximizes aggregate social welfare, and the fixed cost of entry is strictly

positive, then (a) the price cap is below the monopoly price of the lowest cost type and thus induces pooling among all non-excluded types, and (b) some higher cost types must be excluded, provided that not all types would pool at the cap were no exclusion to occur (i.e., provided that the sub-monopoly price that generates zero profit for the highest possible cost type is above the monopoly price of the lowest possible cost type). This proposition demonstrates that exclusion of higher cost types can be optimal in some settings.<sup>6</sup> The setting of linear demand and a uniform distribution is often treated in the literature, and so it may be of some separate interest that we establish conditions for this example under which a price cap is used and higher cost types are excluded.<sup>7</sup>

The described results characterize optimal regulatory policy for market settings in which exclusion for some cost types is feasible. Our results thus directly apply when the monopolist provides an inessential service for a given market or region. Since [Baron and Myerson \(1982\)](#) also focus on settings where exclusion is feasible, our findings also offer a characterization of how their analysis extends when transfers are not feasible.

We are interested as well in the “no-exclusion” scenario, wherein the regulator is constrained to ensure that the monopolist earns non-negative profit while providing positive output under all cost realizations. To motivate this scenario, we note that the monopolist may provide essential services with poor substitution alternatives.<sup>8</sup> An advantage of our three-step approach is that our analysis also can be used to characterize optimal regulatory policies for the no-exclusion scenario. To this end, we may refer to our first proposition for the special case in which the upper-bound cost type equals the highest cost type in the full support. Our first proposition then provides conditions under which optimal regulation for the no-exclusion scenario takes the form of a price-cap policy, where the price cap is set at the second-best level that generates zero profit for a monopolist with the highest possible cost type. Likewise, we can facilitate the application of our results to this scenario by using the three approaches described above. Thus, for example, the described price-cap allocation is optimal for the no-exclusion scenario if the demand function takes a linear, constant elastic-

---

<sup>6</sup>See [Armstrong \(1996\)](#) for an analysis of optimal exclusion in the different context of a model of multi-product nonlinear pricing when the type space is multidimensional.

<sup>7</sup>As discussed below, [Alonso and Matouschek \(2008\)](#) consider this example for a regulation model without transfers but in the absence of an IR constraint. [Baron and Myerson \(1982\)](#) also illustrate their findings using this example.

<sup>8</sup>Other motivations exist as well. For example, the regulator may face a profit constraint whereby the monopolist must enjoy a minimum rate of return for capital expenditures. Alternatively, we can imagine an extended setting in which the legislature chooses a regulatory scheme that is implemented with limited discretion by the regulator, where the regulator faces a commitment problem and cannot induce shut down. For example, suppose that if the legislature selects a price-cap scheme, then the regulator must utilize such a scheme but has some limited discretion concerning the level at which the price cap is placed. If the regulator is unable to commit to allow the monopolist to shut down, then the regulator in effect can only impose a price cap at a level that avoids shut down for any cost type.

ity or log form and if a simple inequality condition holds, respectively, where the inequality condition is sure to hold if the density is non-decreasing over the full support.

Our work is related to research on optimal delegation. The delegation literature begins with [Holmstrom \(1977\)](#), who considers a setting in which a principal faces a privately informed and biased agent and in which contingent transfers are infeasible. The principal then selects a set of permissible actions from the real line, and the agent selects his preferred action from that set after privately observing the state of nature.<sup>9</sup> A key goal in this literature has been to identify general conditions under which the principal optimally defines the permissible set as an interval. [Alonso and Matouschek \(2008\)](#) consider a setting with quadratic utility functions and provide necessary and sufficient conditions for interval delegation to be optimal. Extending the Lagrangian techniques of [Amador et al. \(2006\)](#), [Amador and Bagwell \(2013\)](#) consider a general representation of the delegation problem and establish necessary and sufficient conditions for the optimality of interval delegation.<sup>10</sup>

Our analysis of the regulator’s truncated problem builds on the Lagrangian methods used by [Amador and Bagwell \(2013\)](#), but a novel feature of the current paper is that the analysis is extended to include an ex post participation constraint.<sup>11</sup> Indeed, and as mentioned previously, this constraint would be violated for higher cost types, if we were to apply results from Amador and Bagwell for the delegation problem without such a constraint. A further distinction of the current paper is that, in our analysis of the regulator’s problem, we allow for the possibility of excluded types and show further that actual exclusion can be optimal. In that case, the regulation contract can be understood as providing a disconnected set of quantities, namely, a quantity of zero for excluded types combined with an interval of positive quantities for non-excluded types. The optimal regulation contract is then clearly distinct from an interval allocation. We expect our methods will facilitate the application of optimal

---

<sup>9</sup>A large literature follows Holmstrom’s work. See, for example, [Amador and Bagwell \(2012\)](#), [Amador and Bagwell \(2018\)](#), [Amador et al. \(2018\)](#), [Amador et al. \(2006\)](#), [Ambrus and Egorov \(2017\)](#), [Armstrong and Vickers \(2010\)](#), [Burkett \(2016\)](#), [Frankel \(2014\)](#), [Frankel \(2016\)](#), [Guo \(2016\)](#), [Koessler and Martimort \(2012\)](#), [Martimort and Semenov \(2006\)](#), [Melumad and Shibano \(1991\)](#) and [Mylovanov \(2008\)](#). Related themes also arise in repeated games with private information; see [Athey et al. \(2004\)](#), [Athey et al. \(2005\)](#) and [Halac and Yared \(2019\)](#).

<sup>10</sup>We note that a cap can be understood as a form of interval delegation, in which the maximum (minimum) action is defined by the cap (the lowest “flexible” choice for any agent type).

<sup>11</sup>[Amador and Bagwell \(2018\)](#) also build on the Lagrangian methods used by [Amador and Bagwell \(2013\)](#). The focus of [Amador and Bagwell \(2018\)](#), however, is very different from that of the current paper. [Amador and Bagwell \(2018\)](#) consider an optimal delegation problem with a two-dimensional action set, where one of the actions corresponds to “money burning,” and they provide sufficient conditions under which money burning expenditures are used in an optimal delegation contract. Building on work by [Ambrus and Egorov \(2017\)](#), [Amador and Bagwell \(2018\)](#) also consider an application with an ex ante participation constraint under the assumption that ex ante (non-contingent) transfers are feasible. The participation constraint can then be addressed using standard methods. In the present paper, by contrast, the participation constraint must hold ex post and cannot be addressed using standard methods since transfers are infeasible.



delegation theory to other settings in which participation constraints play an important role.

[Alonso and Matouschek \(2008\)](#) were the first to argue that the monopoly regulation problem can be understood as a delegation problem. As an application of their analysis, they study optimal regulation when costs are privately observed by the regulated firm and transfers are infeasible, and they report conditions under which price-cap regulation is optimal. Our analysis differs in two ways. First, Alonso and Matouschek assume that the monopolist produces regardless of its cost type and do not include a participation constraint. Indeed, their price-cap solution would violate an ex post participation constraint, since the cap is below the marginal cost of the highest-cost firm. By contrast, we include an ex post participation constraint, allow for exclusion, and consider as well the setting in which the ex post participation constraint holds but exclusion is infeasible. When exclusion is not optimal or is infeasible, the optimal price cap in our model is placed at a higher level than in their analysis and generates zero profit for the highest-cost firm. Second, Alonso and Matouschek assume that demand is linear and the regulator maximizes aggregate social surplus. We consider a more general family of demand functions and regulator objectives, and we provide conditions under which exclusion is optimal when demand is linear.

Recent work by [Kolotilin and Zapechelnyuk \(2019\)](#) is also related. They examine optimal delegation in a “linear delegation” framework and, as an application, provide conditions under which a price cap is the optimal regulatory policy in a delegation setting with a participation constraint. The two papers are complementary. We highlight three distinct features of our analysis. First, following [Baron and Myerson \(1982\)](#), we assume that the monopolist has a non-negative fixed cost; by contrast, [Kolotilin and Zapechelnyuk \(2019\)](#) build from the assumption that the monopolist has no fixed costs. Second, the linear delegation framework corresponds in the regulation setting to the family of demand functions that we identify under which the sufficient conditions for our propositions hold if a simple inequality is satisfied; however, as noted above, we can go beyond this family and check the sufficient conditions for our propositions directly, as we do for the exponential demand function. Third, the two papers employ different proof methods: we analyze the delegation problem directly using a Lagrangian approach, whereas [Kolotilin and Zapechelnyuk \(2019\)](#) analyze the delegation problem by drawing a novel link to the literature on Bayesian persuasion.

The remainder of the paper is organized as follows. Section 2 sets up the regulator’s problem, and Section 3 characterizes the optimal regulatory policy when attention is restricted to allocations that can be induced by caps. Section 4 then focuses on the regulator’s truncated problem and develops general sufficient conditions for the optimality of a cap allocation in the set of all allocations that satisfy incentive compatibility and participation constraints. Section 5 considers the global optimality of the cap allocation when attention is widened to



include all types and develops further results and approaches that facilitate the application of our findings. Section 6 identifies conditions under which actual exclusion does or does not occur, respectively. Section 7 concludes. The Appendix contains remaining proofs.

## 2 The Regulator's Problem

In this section, we present our basic model and formally define the problem that confronts the regulator. We also identify the bias in the monopolist's unrestricted output choice.

We consider a monopolist facing an inverse demand function given by  $P(q)$  where  $q$  is the quantity produced. The production quantity  $q$  resides in the set  $Q \equiv [0, q_{max}]$ , which is an interval of the real line with non-empty interior. The function  $P(q)$  is well-defined and finite for all  $q \in (0, q_{max}]$ .

We assume that the monopolist's marginal cost of production is constant and given by  $\gamma$ . The marginal cost  $\gamma$  is private information to the monopolist and is distributed over the support  $\Gamma = [\underline{\gamma}, \bar{\gamma}]$  where  $\bar{\gamma} > \underline{\gamma} > 0$  with a differentiable cumulative distribution function  $F(\gamma)$ . The associated density,  $f(\gamma) \equiv F'(\gamma)$ , is strictly positive and differentiable.

We assume that the regulator has no access to transfers or taxes, and can only impose restrictions on the quantity produced by the monopolist. As discussed in the Introduction, our no-transfers assumption means that the regulator cannot impose taxes or subsidies, and it implicitly implies as well that the monopolist cannot use an access fee. We thus assume that the monopolist selects a uniform price, with the regulator determining the feasible menu of such prices through the selection of a feasible menu of quantities. We allow that the regulator's objective is to maximize a weighted social welfare function in which profits receive weight  $\alpha \in (0, 1]$ . The regulator maximizes aggregate social surplus when  $\alpha = 1$  and gives greater weight to consumer interests when  $\alpha < 1$ .

We impose the following assumptions on primitives:

**Assumption 1.** *We impose the following assumptions:*

- (a)  $P(q)$  is twice-continuously differentiable for  $q \in (0, q_{max}]$  with  $P'(q) < 0 < P(q)$ .
- (b)  $\lim_{q \downarrow 0} P(q) > \bar{\gamma}$  and  $P(q_{max}) < \underline{\gamma}$ .
- (c) There exist functions  $b(q)$ ,  $v(q)$ , and  $w(\gamma, q)$  which are twice-continuously differentiable

for  $q \in (0, q_{max}]$  and that satisfy

$$\begin{aligned} b(q) &\equiv P(q)q, \\ v(q) &\equiv \int_0^q P(z)dz - P(q)q, \\ w(\gamma, q) &\equiv -\gamma q + b(q) + \frac{1}{\alpha}v(q), \end{aligned}$$

with  $\lim_{q \downarrow 0} b(q) = 0$  and  $\lim_{q \downarrow 0} v(q) = 0$ . We define  $b(0) = v(0) = w(\gamma, 0) = 0$ .

(d)  $b''(q) < 0$  and  $w_{qq}(\gamma, q) = b''(q) + \frac{1}{\alpha}v''(q) \leq 0$  for all  $q \in (0, q_{max}]$ .

(e)  $w_q(\gamma, q_{max}) < 0$ .

In this assumption,  $b(q)$  defines the total revenue for the monopolist,  $v(q)$  represents consumer surplus, and  $w(\gamma, q)$  represents the welfare to the regulator (gross of the fixed cost). Using Assumption 1, we obtain that

$$v'(q) = -qP'(q) > 0 \text{ for all } q > 0.$$

Similarly, using Assumption 1, second derivatives take the following forms and signs:

$$\begin{aligned} w_{qq}(\gamma, q) &= b''(q) + \frac{1}{\alpha}v''(q) \\ &= P''(q)q + 2P'(q) - \frac{1}{\alpha}[P''(q)q + P'(q)] \leq 0 \text{ for } q > 0. \end{aligned}$$

Notice that  $P'(q) < 0$  implies that  $w$  is strictly concave when  $\alpha = 1$ . Importantly, we make no assumption as regards the sign of  $v''(q)$ . If marginal revenue is steeper than demand (i.e.,  $b''(q) < P'(q)$ ), then  $v''(q) > 0$ .<sup>12</sup> For example, as we discuss in greater detail below,  $v''(q) > 0$  when demand is linear, and  $v''(q) < 0$  when demand exhibits constant elasticity.

Assumption 1 also includes various regularity conditions. According to part (b), the inverse demand function exceeds the highest marginal cost for quantities that are sufficiently close to zero and falls below the lowest marginal cost for quantities that are sufficiently close to  $q_{max}$ . Part (e) ensures that the welfare-maximizing quantity is below  $q_{max}$ , even when marginal cost is at its lowest possible value.

We envision the regulator as choosing a menu of permissible outputs, with the understanding that a monopolist with cost type  $\gamma$  selects its preferred output from this menu. Thus, if the regulator seeks to assign an output  $q(\gamma)$  to a monopolist with type  $\gamma$ , then an incentive compatibility constraint must be satisfied. As well, if the regulator seeks a positive

---

<sup>12</sup> This condition holds if the demand function is log-concave but fails otherwise.

output from a monopolist with type  $\gamma$ , then type  $\gamma$  must earn more by producing  $q(\gamma) > 0$  than by shutting down and avoiding the fixed cost of production,  $\sigma \geq 0$ .

We allow that the regulator may choose a menu of permissible outputs such that some types produce zero output, incur no fixed cost, and thus earn a profit of zero. That is, the regulator may “exclude” some types from production.

The *regulator’s problem* can then be written as follows:

$$\begin{aligned}
(\text{P1}) \quad & \max_{q: \Gamma \rightarrow Q} \int_{\Gamma} (w(\gamma, q(\gamma)) - \mathbf{1}(q(\gamma))\sigma) dF(\gamma) \quad \text{subject to:} \\
& \gamma \in \arg \max_{\tilde{\gamma} \in \Gamma} -\gamma q(\tilde{\gamma}) + b(q(\tilde{\gamma})) - \mathbf{1}(q(\tilde{\gamma}))\sigma \text{ for all } \gamma \in \Gamma \\
& 0 \leq -\gamma q(\gamma) + b(q(\gamma)) - \mathbf{1}(q(\gamma))\sigma, \text{ for all } \gamma \in \Gamma
\end{aligned}$$

where  $\mathbf{1}(\cdot)$  is an indicator function such that  $\mathbf{1}(q) = 1$  if  $q > 0$  and  $\mathbf{1}(q) = 0$  if  $q = 0$ .

The first constraint in this problem is the incentive compatibility constraint, while the second constraint is the ex post participation or individual rationality (IR) constraint. Notice that the IR constraint requires that if a type produces, it needs to earn enough profit to cover its fixed cost,  $\sigma$ . Notice also that the constraints allow for the possibility of types for which  $q(\gamma) = 0$ , since the IR constraint as represented here holds when  $q(\gamma) = 0$ . We say that an allocation is *feasible* if it satisfies both of these constraints.

**The flexible allocation.** Before moving on to characterize the solution to the regulator’s problem, it is convenient to define  $q_f(\gamma)$  as the allocation that a monopolist would choose if it were forced to produce but were otherwise unrestricted by the regulator. To this end, we let  $\pi(\gamma, q)$  be the monopolist’s profit function (gross of the fixed cost),

$$\pi(\gamma, q) \equiv -\gamma q + b(q),$$

and we then define the monopolist’s flexible allocation as

$$q_f(\gamma) = \arg \max_{q \in Q} \pi(\gamma, q).$$

The flexible allocation is simply the monopoly output as a function of the monopolist’s cost type. The associated first-order condition is given by

$$b'(q) - \gamma = 0.$$

We note that the  $\lim_{q \rightarrow 0} P(q) > \bar{\gamma}$  and  $b(0) = 0$  imply that  $q_f(\bar{\gamma}) > 0$ . Since  $P(q_{max}) < \underline{\gamma}$ ,

we know that  $q_f(\underline{\gamma}) < q_{max}$ . With these boundary results in place, we have that  $q_f(\gamma)$  is differentiable, with  $q'_f(\gamma) = 1/b''(q_f(\gamma)) < 0$  and  $q_f(\gamma) \in (0, q_{max})$  for all  $\gamma \in \Gamma$ . Note as well that  $P(q_f(\gamma)) > \gamma$  and thus  $\pi(\gamma, q_f(\gamma)) = -\gamma q_f(\gamma) + b(q_f(\gamma)) > 0$  for all  $\gamma \in \Gamma$ .

We further assume that it is optimal for all types to produce if given the ability to set their monopolist quantity:

**Assumption 2.** *For all types  $\gamma \in \Gamma$ ,  $\pi(\gamma, q_f(\gamma)) > \sigma$ .*

An implication of Assumption 2 is that, for any given cost type, the regulator's welfare is higher when a monopolist with that cost type sets its monopoly output than when it shuts down and produces zero output. Thus, if the solution to the regulator's problem excludes a given cost type from production, then it must be that the regulator is able to improve the allocation for other cost types through this means.

Given the interiority of  $q_f(\gamma)$ , we may use the associated first-order condition and establish the following relationship:

$$\begin{aligned} w_q(\gamma, q_f(\gamma)) &= \frac{1}{\alpha} v'(q_f(\gamma)) \\ &= -\frac{1}{\alpha} P'(q_f(\gamma)) q_f(\gamma) \\ &= \frac{1}{\alpha} [P(q_f(\gamma)) - \gamma] > 0 \text{ for all } \gamma \in \Gamma. \end{aligned}$$

Thus, the regulator model is characterized by downward or *negative bias*: the agent's (i.e., the monopolist's) preferred  $q$  is too low from the principal's (i.e., the regulator's) perspective.

The presence of negative bias suggests the possibility of a solution that imposes a lower bound on  $q$  for higher types (or equivalently a cap on the price for higher types). But note also that the unrestricted monopolist profits are decreasing in  $\gamma$ ; thus, it is also possible that such a regulatory restriction could exclude higher-cost types from producing, if as a consequence they are unable to cover their fixed cost of production.

We show now that, if any exclusion occurs, then the excluded types are always defined by a threshold type,  $\gamma_t \in \Gamma$ :

**Lemma 1.** *In any feasible allocation  $q(\cdot)$ , there exists a cut-off  $\gamma_t \in [\underline{\gamma}, \bar{\gamma}]$  such that  $q(\gamma) = 0$  for  $\gamma > \gamma_t$  and  $q(\gamma) > 0$  for  $\gamma < \gamma_t$ . In addition, if  $\gamma_t \in (\underline{\gamma}, \bar{\gamma})$ , then  $-\gamma_t q(\gamma_t) + b(\gamma_t) = \sigma$ .*

*Proof.* Suppose to the contrary that for some  $\gamma_1$  and  $\gamma_2$  with  $\underline{\gamma} \leq \gamma_1 < \gamma_2 \leq \bar{\gamma}$ , we have that  $q(\gamma_1) = 0 < q(\gamma_2)$ . A monopolist with type  $\gamma_1$  would then gain by violating the incentive compatibility constraint and selecting instead the output intended for type  $\gamma_2$ :

$$-\gamma_1 q(\gamma_2) + b(q(\gamma_2)) - \sigma > -\gamma_2 q(\gamma_2) + b(q(\gamma_2)) - \sigma \geq 0$$

where the first inequality follows since  $q(\gamma_2) > 0$  and  $\gamma_1 < \gamma_2$ , and the second inequality follows from the IR constraint for a monopolist with type  $\gamma_2$ . As a result, type  $\gamma_1$  prefers to produce  $q(\gamma_2)$  rather than not producing and getting a payoff of 0.

For the second part, suppose that for  $\gamma_t \in (\underline{\gamma}, \bar{\gamma})$ ,  $-\gamma_t q(\gamma_t) + b(q(\gamma_t)) > \sigma$ . Then, for all sufficiently small  $\epsilon > 0$ , we have that  $-(\gamma_t + \epsilon)q(\gamma_t) + b(q(\gamma_t)) > \sigma$ . As a result, type  $\gamma_1 = \gamma_t + \epsilon$  will prefer to produce rather than not, a contradiction of the cut-off property. Suppose instead that  $-\gamma_t q(\gamma_t) + b(q(\gamma_t)) < \sigma$ , so type  $\gamma_t$  strictly prefers not to produce. For all sufficiently small  $\epsilon > 0$ , we have that  $-(\gamma_t - \epsilon)q(\gamma_t - \epsilon) + b(q(\gamma_t - \epsilon)) > \sigma$ , by the cut-off property and strict monotonicity of the profit function in  $\gamma$  when  $q > 0$ . But this implies that  $-\gamma_t q(\gamma_t - \epsilon) + b(q(\gamma_t - \epsilon)) > \sigma - \epsilon q(\gamma_t - \epsilon)$ , and thus, for sufficiently small  $\epsilon$ , type  $\gamma_t$  would strictly prefer to produce given Assumption 2 and choose type  $\gamma_t - \epsilon$ 's choice, a violation of feasibility.  $\square$

If we were to ignore the IR constraint, the regulator's problem would fit into the framework developed by [Amador and Bagwell \(2013\)](#), and we could use the sufficiency theorems in that paper to derive conditions under which a simple cap allocation is optimal.<sup>13</sup> However, as we show below, the IR constraint will always be violated if ignored.

### 3 Optimality Within the Set of Cap Allocations

In this section, we study cap allocations when the IR constraint is ignored and also when exclusion is possible. Our analysis clarifies the role of the IR constraint and identifies a candidate allocation for the solution of the regulator's problem.

#### 3.1 The case without an IR constraint

It is helpful to solve the regulator's problem under the restriction that the regulator can choose only among cap allocations, while ignoring the IR constraint. Let us define a cap allocation as follows:

**Definition 1.** A *cap allocation* indexed by  $x$  is an allocation  $q_c(\gamma; x)$  such that

$$q_c(\gamma; x) = \begin{cases} q_f(\gamma) & \text{if } q_f(\gamma) \geq x \\ x & \text{otherwise} \end{cases}$$

---

<sup>13</sup>One further difference is that the flexible allocation (i.e., the ideal allocation for the monopolist or agent) is upward sloping in the framework of [Amador and Bagwell \(2013\)](#) while the flexible allocation is downward sloping in the current setting. This difference could be easily addressed with a straightforward notational modification, in which  $q$  is re-defined as the extent to which actual output falls short of some upper bound.

for all  $\gamma \in \Gamma$ .

It is straightforward to confirm that a cap allocation is always incentive compatible. Types that would prefer to produce an amount higher than  $x$  are unconstrained and thus choose to produce their flexible output. Types that would prefer to produce an amount lower than  $x$  are constrained and end up producing  $x$ . Naturally, this implies that there exists a critical type  $\gamma_c$ , defined as follows:<sup>14</sup>

**Definition 2.** Given  $x \in Q$ , let  $\gamma_c(x)$  be the unique value in  $\Gamma$  such that  $q_f(\gamma) > x$  for all  $\gamma \in [\underline{\gamma}, \gamma_c(x))$  and  $q_f(\gamma) < x$  for all  $\gamma \in (\gamma_c(x), \bar{\gamma}]$ .

We allow in the definition of  $\gamma_c(x)$  that  $\gamma_c(x) = \underline{\gamma}$ , in which case  $x \geq q_f(\underline{\gamma})$ , so that the flexible output for all types above  $\underline{\gamma}$  is below  $x$ . Notice also that the allocation  $q_c(\gamma; x)$  actually defines a quantity floor rather than a cap. We still refer to this allocation as a cap allocation, since it corresponds to a cap on permissible prices and links thereby with the literature on price-cap regulation. Note also that the cap allocation only has bite in restricting the monopolist's choice if  $x > q_f(\bar{\gamma})$ , as otherwise the monopolist selects its flexible output for all types.

We define an *optimal simple cap allocation* to be an optimal cap allocation when the IR constraint is ignored and all types produce. That is, the optimal simple cap allocation solves

$$\max_{x \geq q_f(\bar{\gamma})} W^c(x)$$

where  $W^c(x)$  represents the regulator's welfare:

$$W^c(x) \equiv \int_{\underline{\gamma}}^{\gamma_c(x)} w(\gamma, q_f(\gamma)) dF(\gamma) + \int_{\gamma_c(x)}^{\bar{\gamma}} w(\gamma, x) dF(\gamma) - \sigma$$

The following lemma provides a necessary condition for an optimal simple cap allocation:<sup>15</sup>

**Lemma 2.** *The cap allocation indexed by  $x$  is an optimal simple cap allocation only if  $x > q_f(\bar{\gamma})$  and*

$$\int_{\gamma_c(x)}^{\bar{\gamma}} w_q(\gamma, x) dF(\gamma) = 0$$

*Proof.* In the appendix. □

---

<sup>14</sup>Here and in the rest of the paper, we use the convention that the intervals  $[x, x)$  and  $(x, x)$  correspond to the empty set.

<sup>15</sup>The existence of an optimal simple cap allocation follows from standard arguments, given Assumption 1. The first-order condition presented in Lemma 2 is necessary but not sufficient for the characterization of an optimal simple cap allocation, since the first-order condition could also characterize a minimum.

In the absence of a participation constraint, we could use results from [Amador and Bagwell \(2013\)](#) and establish a general set of environments under which the optimal simple cap allocation is optimal over the full class of incentive compatible allocations. As we now argue, however, the presence of an IR constraint implies that the optimal simple cap allocation is not feasible.

The basic point can be understood using a graphical argument. The graph on the right in Figure 1 illustrates the optimal simple cap allocation in bold (for the case where  $\gamma_c$  is in the interior of  $\Gamma$ ). This allocation is illustrated relative to the flexible allocation,  $q_f(\gamma)$ , and the regulator's ideal (i.e., efficient) allocation,  $q_e(\gamma)$ , which we define as the allocation that maximizes  $w(\gamma, q)$ .<sup>16</sup> Notice that  $q_e(\gamma)$  is downward sloping and that  $q_e(\gamma) > q_f(\gamma)$ , where the inequality reflects the aforementioned downward bias. For given  $\gamma$ ,  $q_e(\gamma)$  induces a price equal to marginal cost (i.e.,  $P(q_e(\gamma)) = \gamma$ ) when  $\alpha = 1$ . When  $\alpha < 1$ , the regulator's ideal allocation entails even higher quantities and thus drives price below marginal cost. The optimal simple cap allocation is such that the cap is ideal for the regulator on average for affected types (i.e., for  $\gamma \geq \gamma_c$ ). The graph on the left in Figure 1 illustrates the same information in terms of the induced prices, which are also depicted in bold. As this graph illustrates, the optimal simple cap allocation places the price cap at a level that is ideal for the principal on average for affected types. This graph also suggests that the participation constraint is violated for the highest types when the optimal simple cap allocation is used. For type  $\bar{\gamma}$ , the optimal price cap lies below the regulator's ideal price,  $P(q_e(\bar{\gamma}))$ , which equals  $\bar{\gamma}$  when  $\alpha = 1$  and is less than  $\bar{\gamma}$  when  $\alpha < 1$ . The optimal price cap is thus strictly below  $\bar{\gamma}$ ; hence, since the fixed cost  $\sigma$  is non-negative, the IR constraint must fail for the highest-cost type when the optimal simple cap allocation is used.

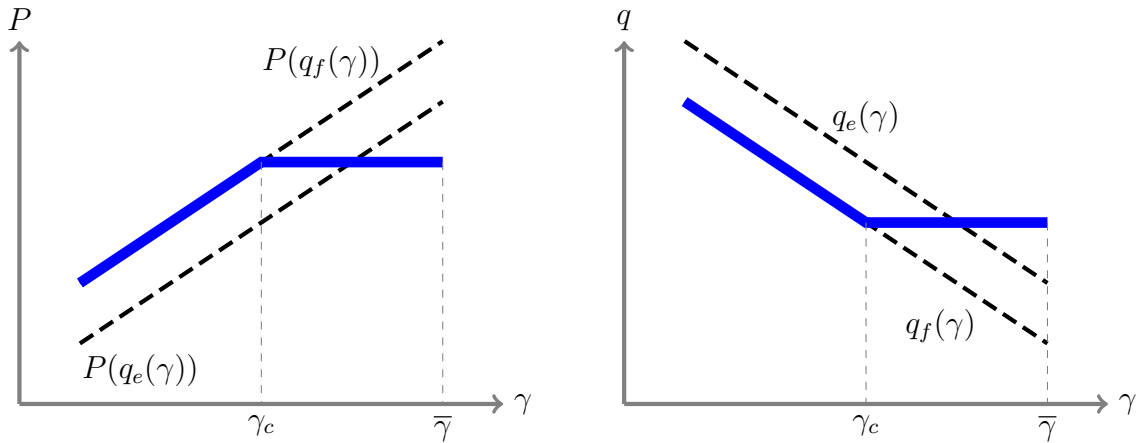


Figure 1: Optimal Simple Cap Allocation Fails IR.

<sup>16</sup>We assume for this graphical analysis that  $q_e(\gamma)$  is uniquely determined.



To develop this point with full details, let  $\pi_c(\gamma; x) = b(q_c(\gamma; x)) - \gamma q_c(\gamma; x)$ . The participation constraint is equivalent to  $\pi_c(\gamma; x) \geq \sigma$  for all  $\gamma \in \Gamma$ . Note that  $\pi_c(\gamma; x)$  is continuous, and that

$$\pi'_c(\gamma; x) = \begin{cases} (b'(q_c(\gamma; x)) - \gamma)q'_c(\gamma; x) - q_c(\gamma; x) = -q_f(\gamma) < 0 & ; \gamma \in (\underline{\gamma}, \gamma_c(x)) \\ -q_f(\gamma_c(x)) < 0 & ; \gamma \in (\gamma_c(x), \bar{\gamma}) \end{cases}$$

Hence,  $\pi_c$  is strictly decreasing in  $\gamma$ . So to check whether the IR constraint holds it suffices to check whether  $\pi_c(\bar{\gamma}; x) \geq \sigma$ , that is, whether the allocation is individually rational for the highest cost type. We have the following lemma:

**Lemma 3.** *The optimal simple cap allocation indexed by  $x$  violates the IR constraint for the highest types.*

*Proof.* Recall from Lemma 2 that  $x > q_f(\bar{\gamma})$  and  $\int_{\gamma_c(x)}^{\bar{\gamma}} w_q(\gamma, x) dF(\gamma) = 0$ . Using  $w_{q\gamma}(\gamma, q) = -1 < 0$ ,  $w_q(\bar{\gamma}, x) < 0$  follows. Next, observe that  $w_q(\bar{\gamma}, x) = -\bar{\gamma} + P(x) + (\frac{1-\alpha}{\alpha})(-P'(x)x) < 0$ , and thus  $P(x) - \bar{\gamma} < (\frac{1-\alpha}{\alpha})(P'(x)x) \leq 0$  given that  $P'(x) < 0$  and  $\alpha \in (0, 1]$ . We then have that  $\pi(\bar{\gamma}, x) = (P(x) - \bar{\gamma})x < 0 \leq \sigma$ , and thus the IR constraint is violated for the highest type.  $\square$

There are two ways a regulator could in principle deal with the problem that the optimal simple cap allocation violates the IR constraint. First, it could decide not to be so tough, and choose a cap that gives sufficient flexibility so that all types choose to produce. Alternatively, it could choose a cap that is sufficiently tight that some types choose not to produce.<sup>17</sup> This leads us to consider the “best” cap allocation that satisfies the IR constraint while allowing types to be excluded from production. We proceed to characterize the class of allocations with caps and exclusion.

### 3.2 IR constraint and exclusion

Consider a situation where the regulator chooses a cap on the price that can be charged, and as a result, some high-cost types may choose not to produce. This is a *cap allocation with potential exclusion*, and it is defined by a quantity  $x$  such that any type is free to choose between producing a quantity higher or equal to  $x$ , or not producing at all:

---

<sup>17</sup>As mentioned in the Introduction, we are also interested in the scenario where the regulator is constrained to ensure that all types produce. An advantage of our solution approach is that the optimal form of regulation for the “no-exclusion” scenario can be characterized as a by-product of our proof for the general case in which exclusion is allowed.

**Definition 3.** A cap allocation with potential exclusion indexed by  $x$  is an allocation  $q(\gamma; x)$  such that

$$q(\gamma; x) = \begin{cases} q_f(\gamma) & ; \text{if } q_f(\gamma) \geq x, \\ x & ; \text{if } q_f(\gamma) < x \text{ and } -\gamma x + b(x) - \sigma \geq 0, \\ 0 & ; \text{otherwise,} \end{cases}$$

for all  $\gamma \in \Gamma$ .

Note that similar to before, a cap allocation with potential exclusion is incentive compatible. Without loss of generality, we can restrict attention to cap allocations such that  $x \geq \underline{q} \equiv q_f(\bar{\gamma})$ , as no type will ever choose to produce below  $q_f(\bar{\gamma})$ , if given the choice to produce more. Similarly, we can restrict attention to cap allocations such that  $x \leq \bar{q}$  where  $\bar{q} > q_f(\underline{\gamma})$  is the value that satisfies  $-\underline{\gamma}\bar{q} + b(\bar{q}) = \sigma$ . Imposing a bound  $x$  above  $\bar{q}$  is equivalent to assigning no production for all types (as not even the lowest cost type is willing to produce that much), and hence considering restrictions above that is unnecessary. Note that our assumptions guarantee  $\bar{q} \in Q$ .

Figure 2 presents a graphical representation of a cap allocation with exclusion where a non-zero measure of types are excluded, some types are constrained at the cap, and some other types are choosing their monopoly allocation. To describe such an allocation, recall that, from Lemma 1, we know that any allocation with exclusion satisfies a threshold property: types above some type  $\gamma_t$  are excluded from production, while types below  $\gamma_t$  produce. Thus, given a bound  $x$ , let  $\gamma_t(x) \in [\underline{\gamma}, \bar{\gamma}]$  be the associated exclusion threshold. That is,  $\gamma_t(x)$  is such that  $\max_{q \geq x} \{-\gamma q + b(q) - \sigma\} < 0$  for all  $\gamma \in (\gamma_t(x), \bar{\gamma}]$  and  $\max_{q \geq x} \{-\gamma q + b(q) - \sigma\} > 0$  for all  $\gamma \in [\underline{\gamma}, \gamma_t(x))$ .

However, not all the types that produce are able to do so at their monopoly level. Types with a cost smaller than  $\gamma_c(x)$  would choose their monopoly level if forced to produce, while types above  $\gamma_c(x)$  would choose the cap if forced to produce. Note that  $\gamma_c(x) \leq \gamma_t(x)$  with strict inequality if  $\underline{q} < x < \bar{q}$ .

We can thus write the welfare generated by a cap allocation with exclusion as:

$$W(x) \equiv \int_{\underline{\gamma}}^{\gamma_c(x)} [w(\gamma, q_f(\gamma)) - \sigma] dF(\gamma) + \int_{\gamma_c(x)}^{\gamma_t(x)} [w(\gamma, x) - \sigma] dF(\gamma) \quad (1)$$

where the first term represents the regulator's payoff from giving flexibility to types below  $\gamma_c(x)$ , the second term represents the payoffs generated from types that produce at the cap,  $x$ , and where the payoff of the excluded types is zero.

Let  $x^*$  be such that  $x^* \in \operatorname{argmax}_{x \in [\underline{q}, \bar{q}]} W(x)$ ; that is,  $x^*$  represents the optimal cap that

could be imposed.<sup>18</sup> Given this cap  $x^*$ , the associated cap allocation  $q^*$  can be written as:

$$q^*(\gamma) = \begin{cases} q_f(\gamma) & \gamma \in [\underline{\gamma}, \gamma_c(x^*)) \\ x^* & \gamma \in [\gamma_c(x^*), \gamma_t(x^*)] \\ 0 & \gamma \in (\gamma_t(x^*), \bar{\gamma}] \end{cases} \quad (2)$$

This cap allocation with exclusion  $q^*$  is our candidate allocation for the solution to the regulator's problem. Our goal is thus to determine sufficient conditions under which we can be certain that  $q^*$  is also optimal within the set of all feasible allocations.

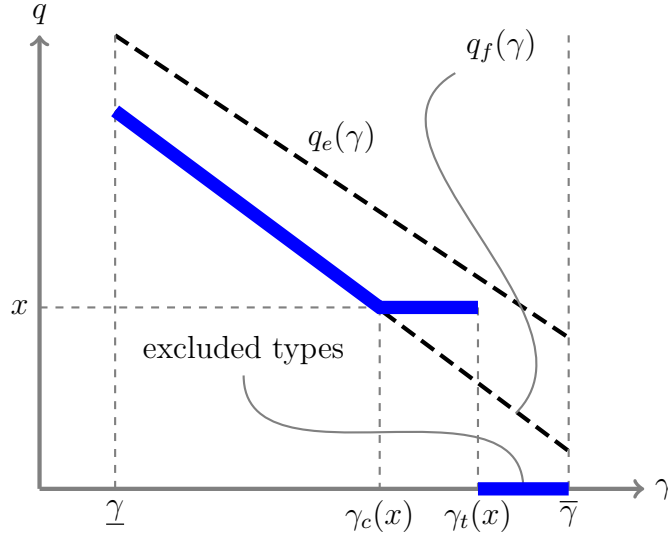


Figure 2: A Cap Allocation with Exclusion. The solid thick line represents a cap allocation with exclusion.

## 4 Towards Sufficient Conditions

We return now to consider the solution to the regulator's problem, Problem P1. As a general matter, we do not know whether a cap allocation with or without exclusion is optimal. Indeed, solving the regulator's problem, Problem P1, directly seems difficult, since the possibility of excluding some types must be considered. We pursue an alternative approach, one that divides the problems into several sub-problems.

The main idea is as follows:

1. Rather than working with the lower bound on production, we work with the excluded types directly. In particular, based on Lemma 1, we fix a given threshold for excluded

---

<sup>18</sup>The existence of  $x^*$  follows from standard arguments, given Assumption 1.

types,  $\gamma_t$ , and consider the problem of allocating production for types below  $\gamma_t$  while ignoring the allocation for types higher than  $\gamma_t$ . That is, we restrict attention to the set of types  $[\underline{\gamma}, \gamma_t]$  and study the problem for a regulator that only considers types in that truncated set and is not allowed to exclude any types in that set from production. We obtain conditions under which an optimal allocation in this truncated problem is a cap allocation with a cap at a level such that the threshold or upper-bound type,  $\gamma_t$ , is indifferent between producing or not.

2. Next, we argue that such truncated allocation is incentive compatible when extended to the entire set  $[\underline{\gamma}, \bar{\gamma}]$  by giving types above  $\gamma_t$  zero output. This implies that the optimal allocation that results from considering only the truncated set is also optimal when considering the entire set of types for a given level of exclusion.
3. We then look for the best allocation by varying the level of exclusion, which in our case is indexed by  $\gamma_t$ . This is a single variable optimization problem.

Towards this goal, let us first consider the regulator's truncated problem.

## 4.1 The Regulator's Truncated Problem

For this problem, we fix  $\gamma_t \in (\underline{\gamma}, \bar{\gamma}]$  and define  $\Gamma_t(\gamma_t) \equiv [\underline{\gamma}, \gamma_t]$ .<sup>19</sup> The regulator's truncated problem is to find an allocation,  $q_t : \Gamma_t(\gamma_t) \rightarrow Q$ , that maximizes its payoff subject to the feasibility constraints and that no type in set  $\Gamma_t(\gamma_t)$  is excluded.<sup>20</sup> The problem is

$$\begin{aligned} \max_{q_t : \Gamma_t(\gamma_t) \rightarrow Q} \int_{\Gamma_t(\gamma_t)} (w(\gamma, q_t(\gamma)) - \sigma) dF(\gamma) \quad \text{subject to:} \quad (P_t) \\ \gamma \in \arg \max_{\tilde{\gamma} \in \Gamma_t(\gamma_t)} \{-\gamma q_t(\tilde{\gamma}) + b(q_t(\tilde{\gamma})) - \sigma\} \text{ for all } \gamma \in \Gamma_t(\gamma_t) \\ 0 \leq -\gamma q_t(\gamma) + b(q_t(\gamma)) - \sigma, \text{ for all } \gamma \in \Gamma_t(\gamma_t) \end{aligned}$$

Note that differently from Problem P1, in this truncated regulator problem all types are producing – this explains why the indicator functions do not appear in Problem  $P_t$ . Similarly to Subsection 3.1, if we were to look for a simple cap allocation in this truncated problem, the optimal one will violate the IR constraint for the highest cost type, in this case the threshold or upper-bound type,  $\gamma_t$ .

<sup>19</sup>We ignore the case where  $\gamma_t = \underline{\gamma}$ , as this implies that almost all types are excluded, a situation that cannot be optimal under our assumptions.

<sup>20</sup>Note that when looking within the set of cap allocations, it is sufficient to look for a quantity floor in  $[q, \bar{q}]$ . When checking for optimality more generally, we do not impose that restriction, and hence  $q_t : \Gamma(\gamma_t) \rightarrow Q$ .

We conjecture however that a cap allocation where type  $\gamma_t$  is indifferent between producing or not is optimal. Let  $q_i(\gamma_t)$  be the unique value such that  $-\gamma_t q_i(\gamma_t) + b(q_i(\gamma_t)) = \sigma$  and  $q_i(\gamma_t) > q_f(\gamma_t)$ . Thus,  $q_i(\gamma_t)$  is the output level that exceeds  $\gamma_t$ 's monopoly level and ensures that this type is indifferent between producing at that level or not. In other words, it corresponds to a price that equals the average cost for type  $\gamma_t$ . Note that under our assumptions, such  $q_i(\gamma_t) \in Q$  exists.

We define  $\gamma_H(\gamma_t) \in [\underline{\gamma}, \gamma_t]$  to be the value such that  $q_i(\gamma_t) \leq q_f(\gamma)$  for  $\gamma < \gamma_H(\gamma_t)$ , and  $q_i(\gamma_t) \geq q_f(\gamma)$  for  $\gamma > \gamma_H(\gamma_t)$ . Note that  $\gamma_H(\gamma_t) = \gamma_c(q_i(\gamma_t))$  and that  $\gamma_H(\gamma_t) < \gamma_t$  given  $\gamma_t > \underline{\gamma}$ .

With these objects, we can define the *truncated cap allocation*,  $q_t^*(\gamma|\gamma_t)$ :

$$q_t^*(\gamma|\gamma_t) = \begin{cases} q_f(\gamma) & \gamma \in [\underline{\gamma}, \gamma_H(\gamma_t)) \\ q_i(\gamma_t) & \gamma \in [\gamma_H(\gamma_t), \gamma_t] \end{cases} \quad (3)$$

This allocation  $q_t^*(\gamma|\gamma_t)$  is continuous in  $\gamma$  and may feature full pooling of types if  $\gamma_H(\gamma_t) = \underline{\gamma}$ . Note that if  $\gamma_H(\gamma_t)$  is interior to the interval  $\Gamma_t(\gamma_t)$ , then  $q_i(\gamma_t)$  coincides with the flexible quantity chosen by type  $\gamma_H(\gamma_t)$ . Figure 3 displays the two possible cases for  $q_t^*$  for two different values of  $\gamma_t$ . Panel (a) shows the case with partial pooling. Panel (b) shows the cases where  $\gamma_t$  is sufficiently small that full pooling of all types at the cap results.

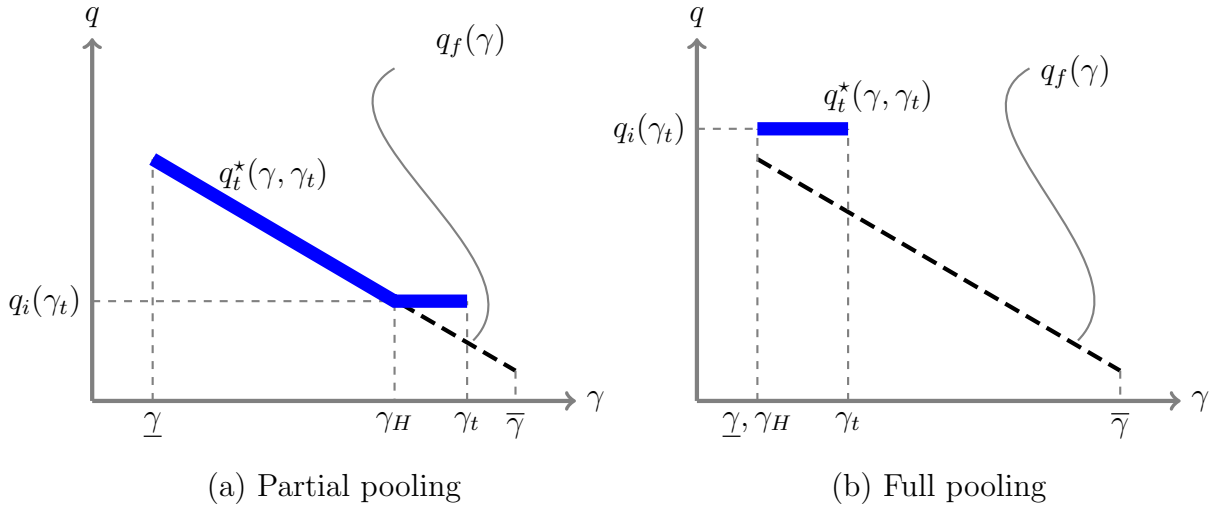


Figure 3: The truncated cap allocation,  $q_t^*(\gamma|\gamma_t)$ .

We would like to find conditions under which  $q_t^*(\gamma|\gamma_t)$  is the optimal solution to the regulator's truncated problem. To present our next result, we require a couple of definitions.

Let

$$G(\gamma|\gamma_t) \equiv -\kappa F(\gamma_t) + \kappa \left[ \frac{\gamma - b'(q_i(\gamma_t))}{\gamma - \gamma_H(\gamma_t)} \right] F(\gamma) + \frac{1}{\gamma - \gamma_H(\gamma_t)} \int_{\gamma_H(\gamma_t)}^{\gamma} w_q(\tilde{\gamma}, q_i(\gamma_t)) f(\tilde{\gamma}) d\tilde{\gamma}, \quad (4)$$

for  $\gamma > \gamma_H(\gamma_t)$  and where, following [Amador and Bagwell \(2013\)](#),  $\kappa$  is a relative concavity parameter defined as

$$\kappa \equiv \min_{q \in Q} \left\{ 1 + \frac{v''(q)}{\alpha b''(q)} \right\}.$$

We let  $G(\gamma_H(\gamma_t)|\gamma_t) \equiv \lim_{\gamma \downarrow \gamma_H(\gamma_t)} G(\gamma|\gamma_t)$ , which exists and is a finite number.

We may now state our general sufficiency result as follows:

**Proposition 1.** (*Sufficient Conditions*) *If*

(i)  $G(\gamma|\gamma_t) \leq G(\gamma_t|\gamma_t)$  for all  $\gamma \in [\gamma_H(\gamma_t), \gamma_t]$ , where  $G$  as given by (4); and

(ii)  $M_1(\gamma) \equiv \kappa F(\gamma) + w_q(\gamma, q_f(\gamma))f(\gamma)$  is non-decreasing in  $\gamma$  for  $\gamma \in [\underline{\gamma}, \gamma_H(\gamma_t))$ ,

then the cap allocation  $q_t^*(\gamma|\gamma_t)$  solves the regulator's truncated problem, Problem  $P_t$ .

*Proof.* In the appendix. □

Our proof approach follows a guess-and-verify structure. To begin, we follow standard methods and re-write the incentive constraint in the regulator's truncated problem as an integral equation and a monotonicity requirement (namely, that  $q_t(\gamma)$  must be non-increasing).<sup>21</sup> Next, we embed the monotonicity requirement into the choice set, and we express the integral equation equivalently in terms of two inequality conditions. The regulator's truncated problem is thereby represented as a maximization problem over functions belonging to a choice set of non-decreasing functions that satisfy three inequality constraints, where one of the constraints is the IR constraint. With the problem set up in this fashion, we conjecture that the cap allocation  $q_t^*(\gamma|\gamma_t)$  is the solution. To confirm this conjecture, we construct multiplier functions for each of the three inequality constraints. Under the conditions stated in Proposition 1 and for the constructed multiplier functions, we find that the multiplier functions are non-decreasing, the corresponding Lagrangian is concave, and the cap allocation satisfies first-order conditions and a complementary slackness condition. Building on work by [Amador and Bagwell \(2013\)](#), we conclude the proof by showing that these findings are sufficient to conclude that  $q_t^*(\gamma|\gamma_t)$  solves the regulator's truncated problem.<sup>22</sup>

<sup>21</sup>We emphasize that feasible allocations may be discontinuous. As illustrated in the intuition developed just below, our proof approach thus must establish that the cap allocation  $q_t^*(\gamma|\gamma_t)$  is optimal among a set of monotone and possibly discontinuous functions.

<sup>22</sup>It is instructive here to compare our regulator's truncated problem, in which transfers are unavailable,

## 4.2 Intuition

We now develop some intuition for the interpretation of Proposition 1. We begin with part (ii). Observe that part (ii) is more easily satisfied when  $\kappa$  is big. Referring to the definition of  $\kappa$ , we thus conclude that part (ii) is more easily satisfied when the minimum value for  $1 + \frac{v''(q)}{\alpha b''(q)}$  is big. Additionally, since  $w_q(\gamma, q_f(\gamma)) > 0$ , we see that part (ii) is also more easily satisfied when the density is non-decreasing for  $\gamma \in [\underline{\gamma}, \gamma_H(\gamma_t))$ .

To see why the relative magnitudes of  $\frac{1}{\alpha}v''(q)$  and  $b''(q)$  and the slope of the density matters, we consider alternatives to the truncated cap allocation. If the truncated cap allocation is to be optimal among all feasible allocations for the regulator's truncated problem, then in particular it must be preferred by the regulator to alternative feasible allocations that are generated by “drilling holes” in the flexible part of the allocation. Figure 4 illustrates one such alternative allocation, in which output levels between  $q_1 \equiv q_f(\gamma_1)$  and  $q_2 \equiv q_f(\gamma_2)$  are prohibited and where  $\underline{\gamma} < \gamma_1 < \gamma_2 < \gamma_H$ . There then exists a unique type  $\tilde{\gamma} \in (\gamma_1, \gamma_2)$  that is indifferent between  $q_1$  and  $q_2$ . The alternative allocation thus induces a “step” at  $\tilde{\gamma}$ , with the allocation  $q_1$  selected by  $\gamma \in [\gamma_1, \tilde{\gamma})$  and the allocation  $q_2$  selected by  $\gamma \in [\tilde{\gamma}, \gamma_2]$ , where for simplicity we place type  $\tilde{\gamma}$  with the higher types.

In comparison to the truncated cap allocation, the alternative allocation has advantages and disadvantages. First, the alternative allocation generates output choices for  $\gamma \in [\gamma_1, \tilde{\gamma})$  that are closer to the the regulator's ideal choices for such types; however, the alternative allocation also results in output choices for  $\gamma \in [\tilde{\gamma}, \gamma_2]$  that are further from the regulator's ideal choices for such types. In line with our discussion above, these observations suggest that a non-decreasing density should work in favor of the truncated cap allocation, since the disadvantageous features of the alternative allocation then receive greater probability weight in the regulator's expected welfare. Second, the alternative allocation increases the variance of the induced allocation around  $q_f(\gamma)$  over the interval  $[\gamma_1, \gamma_2]$ . Consistent with our preceding discussion, this effect brings into consideration the relative magnitudes of  $\frac{1}{\alpha}v''(q)$  and  $b''(q)$ , where the latter determines the slope of  $q_f(\gamma)$ . In particular, if  $v(q)$  is concave, then the variance effect should work in favor of the truncated cap allocation, since the regulator would then not welcome an increase in variance. If instead  $v(q)$  is convex, then

---

with the standard (Baron-Myerson) framework in which transfers are available. In the solution approach for the standard framework, the integral equation is substituted into the objective, the IR constraint is shown to bind for the highest type, the IR constraint for the highest type is substituted into the objective, and the resulting objective is then maximized point-wise. If the solution satisfies the monotonicity constraint, then the problem is solved. By contrast, in our no-transfers setting, we cannot substitute the integral equation into the objective, since we do not have a remaining transfer instrument with which to ensure that the solution of the resulting optimization problem satisfies the integral equation. For the same reason, we cannot substitute the IR constraint for the highest type into the objective. Indeed, as a general matter, when transfers are unavailable it is no longer obvious that the IR constraint for the highest type must bind.



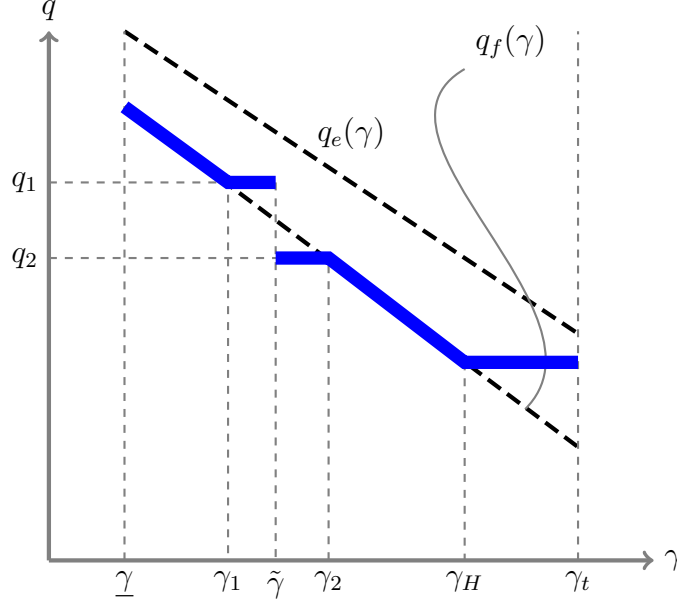


Figure 4: Drilling a hole (with  $\sigma > 0$ ).

the regulator would benefit from the greater variance afforded by the alternative allocation, with the benefit to the regulator being larger when  $\alpha$  is smaller. Based on this perspective, we may understand that the truncated allocation could remain optimal when  $v(q)$  is convex, provided that the density rises fast enough,  $\alpha$  is sufficiently large and/or  $b''(q)$  is large in absolute value (so that  $q_f(\gamma)$  is flat, in which case steps add little variation).

The intuitive discussion presented here considers only a subset of feasible alternative allocations that introduce variations in the flexible region. In our no-transfer setting, the incentive compatibility constraint implies that an allocation must be given by the flexible allocation over any interval for which the allocation is continuous and strictly decreasing; however, an incentive compatible allocation may include many points of discontinuity (steps), where any such point hurdles the flexible allocation as illustrated in Figure 4.<sup>23</sup> Our discussion above considers only an alternative allocation with a single step, but this discussion provides an intuitive foundation for understanding more generally the key forces at play.

We turn now to consider the intuition associated with part (i) of Proposition 1. For type  $\gamma_t$ , the IR constraint holds with equality at the output choices  $q_i(\gamma_t)$  and  $q'$ , where  $q' < q_f(\gamma_t)$  is defined so that type  $\gamma_t$  is indifferent between  $q_i(\gamma_t)$  and  $q'$ ; thus, the IR constraint for type  $\gamma_t$  is satisfied provided that the allocation for this type resides in the interval  $[q', q_i(\gamma_t)]$ . As a general matter, it is not obvious that the IR constraint must bind for type  $\gamma_t$  in our no-transfer setting, since the allocation for this type may be positioned so as to favorably

<sup>23</sup>For further discussion, see [Melumad and Shibano \(1991\)](#).

affect allocations for lower types. Part (i) of Proposition 1 provides conditions under which the solution to the regulator's truncated problem is such that type  $\gamma_t$  selects  $q_i(\gamma_t)$  and has a binding IR constraint.

At a more formal level, we show in the proof that the value of the multiplier function for the IR constraint of type  $\gamma_t$  in fact equals  $G(\gamma_t|\gamma_t)$ . Since the constructed multiplier functions must be non-negative, we thereby confirm that  $G(\gamma_t|\gamma_t) \geq 0$ , with the corresponding interpretation that the shadow price of relaxing the IR constraint for type  $\gamma_t$  is non-negative. Part (i) of Proposition 1 goes further and requires that  $G(\gamma|\gamma_t) \leq G(\gamma_t|\gamma_t)$  for all  $\gamma \in [\gamma_H(\gamma_t), \gamma_t]$ . As confirmed in the proof, this condition ensures that the regulator cannot improve on the cap allocation  $q_t^*(\gamma|\gamma_t)$  by altering the allocation for types  $\gamma \in [\gamma_H(\gamma_t), \gamma_t]$  while respecting the monotonicity requirement.

## 5 Global Optimality

The results of the previous section offer a characterization of the optimal solution given an exogenous amount of exclusion as defined by the fixed threshold or upper-bound type,  $\gamma_t$ . In particular, for every exclusion threshold  $\gamma_t$ , we have found sufficient conditions for the associated truncated cap allocation  $q_t^*$ , defined in (3), to be optimal when restricting attention only to those types not excluded from production. However, it is straightforward to argue now that, given an amount of exclusion, the truncated cap allocation is optimal when attention is widened to include all types. Note that the only potential issue is incentive compatibility. The  $q_t^*$  allocation when extended for all types must remain incentive compatible. But this is straightforward: since type  $\gamma_t$  is indifferent between producing or not, all types above  $\gamma_t$  strictly prefer not to produce, as they face a higher marginal cost.

We have the following result:

**Proposition 2.** *Assume that parts (i) and (ii) of Proposition 1 hold for all  $\gamma_t \in (\underline{\gamma}, \bar{\gamma}]$ . Then the cap allocation with exclusion  $q^*$  defined in (2) solves the regulator problem, Problem P1.*

*Proof.* We know from Lemma 1 that any level of exclusion is given by a threshold  $\gamma_t \in (\underline{\gamma}, \bar{\gamma}]$ . Given any level of exclusion  $\gamma_t$ , the allocation  $q^*(\gamma|\gamma_t)$  defined in equation (3) remains a feasible allocation when the allocation is extended to entire type space by assigning no production to types strictly above  $\gamma_t$ . This follows because type  $\gamma_t$  is indifferent between producing or not in the  $q^*(\gamma|\gamma_t)$  allocation, and thus, all types higher than  $\gamma_t$  strictly prefer not to produce, as prescribed by the allocation.

Thus, for a given level of exclusion,  $\gamma_t$ , Proposition 1 guarantees that the allocation  $q^*(\gamma|\gamma_t)$  extended over the entire type space is optimal within all feasible allocations that

deliver the same level of exclusion.

Note that the allocation  $q^*(\gamma)$  is the optimal among all the  $q^*(\gamma|\gamma_t)$  allocations for all  $\gamma_t \in (\underline{\gamma}, \bar{\gamma}]$ . We can ignore any allocation where  $\gamma_t = \underline{\gamma}$  (that is, full exclusion) as such an allocation is dominated by the fully flexible allocation. As a result  $q^*(\gamma)$  is optimal among the set of all feasible allocations.  $\square$

The following corollary provides easier-to-check conditions for Proposition 1 and 2

**Corollary 1.** *Suppose that  $\kappa \geq 1/2$ . For given  $\gamma_t$ , if  $f(\gamma)$  is non-decreasing for all  $\gamma \in [\underline{\gamma}, \gamma_t]$ , then conditions (i) and (ii) of Proposition 1 hold. If  $f(\gamma)$  is non-decreasing for all  $\gamma \in [\underline{\gamma}, \bar{\gamma}]$ , then the cap allocation with exclusion  $q^*$  is optimal within the set of all feasible allocations.*

*Proof.* In the appendix.  $\square$

## 5.1 A Demand Family

In the Appendix proof of Corollary 1, we show that if the following  $M_2(\gamma)$  function,

$$M_2(\gamma) \equiv \kappa F(\gamma) + \frac{1}{\alpha} v'(q_i(\gamma_t)) f(\gamma) + (\kappa - 1)(\gamma - b'(q_i(\gamma_t))) f(\gamma),$$

is non-decreasing in  $[\gamma_H(\gamma_t), \gamma_t]$ , then part (i) of Proposition 1 holds. We now show that for a demand family (that includes several commonly used examples as we show below),  $M_1(\gamma) = M_2(\gamma)$ ; and thus, if part (ii) of Proposition 1 holds globally for all  $\gamma \in [\underline{\gamma}, \bar{\gamma}]$ , then part (i) holds as well.

Toward this end, we consider a family of demand functions such that

$$\frac{P'(q)}{P(q)} q = a_0 + \frac{b_0}{P(q)} \quad \text{for all } q \in (0, q_{max}] \quad (5)$$

with  $a_0 \neq -1$ . We have the following result:

**Lemma 4.** *Suppose that (5) holds. Then*

$$(a) \quad v(q) = -\frac{a_0}{1+a_0} b(q) - \frac{b_0}{1+a_0} q \text{ for all } q \in Q,$$

$$(b) \quad \kappa = 1 + \frac{1}{\alpha} \frac{v''(q)}{b''(q)} = 1 - \frac{1}{\alpha} \frac{a_0}{1+a_0},$$

$$(c) \quad M_1(\gamma) = M_2(\gamma) \text{ for all } \gamma \in [\underline{\gamma}, \bar{\gamma}].$$

*Proof.* For parts (a) and (b), recall that  $v'(q) = -P'(q)q$  and that  $b'(q) = P(q) + qP'(q)$ . Using equation (5), it follows that, for all  $q \in (0, q_{max}]$ ,

$$\begin{aligned} -v'(q) &= a_0(b'(q) + v'(q)) + b_0 \\ v'(q) &= -\frac{a_0}{1+a_0}b'(q) - \frac{b_0}{1+a_0} \end{aligned} \quad (6)$$

Integrating the above in  $[q_0, q]$  where  $q > q_0 > 0$ , we have that

$$v(q) + \frac{a_0}{1+a_0}b(q) + \frac{b_0}{1+a_0}q = v(q_0) + \frac{a_0}{1+a_0}b(q_0) + \frac{b_0}{1+a_0}q_0$$

From Assumption 1, using the limit condition  $\lim_{q_0 \downarrow 0} v(q_0) = \lim_{q_0 \downarrow 0} b(q_0) = 0$ , we get part (a) for all  $q \in Q$ .

Differentiating (6), we get that  $\frac{1}{\alpha} \frac{v''(q)}{b''(q)} = -\frac{1}{\alpha} \frac{a_0}{1+a_0}$ , and thus part (b) follows.

To show part (c), note that

$$\begin{aligned} w_q(\gamma, q_f(\gamma)) &= \frac{1}{\alpha} v'(q_f(\gamma)) \\ &= -\frac{1}{\alpha} \frac{a_0}{1+a_0} b'(q_f(\gamma)) - \frac{1}{\alpha} \frac{b_0}{1+a_0} \\ &= -\frac{1}{\alpha} \frac{a_0}{1+a_0} \underbrace{(b'(q_f(\gamma)) - b'(q_i))}_{\gamma} - \frac{1}{\alpha} \underbrace{\left[ \frac{a_0}{1+a_0} b'(q_i) - \frac{b_0}{1+a_0} \right]}_{-v'(q_i)} \\ &= (\kappa - 1)(\gamma - b'(q_i)) + \frac{1}{\alpha} v'(q_i) \end{aligned}$$

It follows then that

$$M_1(\gamma) = \kappa F(\gamma) + w_q(\gamma, q_f(\gamma))f(\gamma) = \kappa F(\gamma) + (\kappa - 1)(\gamma - b'(q_i))f(\gamma) + \frac{1}{\alpha} v'(q_i)f(\gamma) = M_2(\gamma)$$

which with  $q_i = q_i(\gamma_t)$  delivers part (c).  $\square$

For the demand family stated in equation (5), we can obtain a general sufficient condition for the results in Propositions 1 and 2 to hold.

**Corollary 2.** *Suppose that  $P$  satisfies (5) for  $a_0 \neq -1$ . If*

$$(2\kappa - 1)f(\gamma) + \frac{1}{\alpha} v'(q_f(\gamma))f'(\gamma) \geq 0 \quad (7)$$

*holds for all  $\gamma \in [\underline{\gamma}, \bar{\gamma}]$ , then conditions (i) and (ii) of Propositions 1 and 2 hold.*

*Proof.* Under this family, we already know that if part (ii) of Proposition 1 holds globally,

then so does part (i). Taking a derivative of  $M_1(\gamma)$  with respect to  $\gamma$ , and using that  $w_q(\gamma, q_f(\gamma)) = \frac{1}{\alpha} v'(q_f(\gamma))$  together with

$$\frac{dv'(q_f(\gamma))}{d\gamma} = \frac{v''(q_f(\gamma))}{b''(q_f(\gamma))}$$

delivers the result.  $\square$

This demand family incorporates several common examples as special cases:

**Linear demand.** Consider  $P(q) = \mu - \beta q$  with  $\mu > \bar{\gamma}$ ,  $\beta > 0$  and  $Q = [0, \mu/\beta - \epsilon]$  for  $\epsilon > 0$  small. For this example,  $q_f(\gamma) = (\mu - \gamma)/(2\beta)$ ,  $v(q) = \beta q^2/2$  and  $\kappa = 1 - \frac{1}{2\alpha}$ . Assumption 1 is satisfied for  $\epsilon > 0$  sufficiently small if  $\alpha \in [\mu/(\mu + \underline{\gamma}), 1]$  where  $1 > \mu/(\mu + \underline{\gamma}) > 1/2$  follows from  $\mu > \underline{\gamma} > 0$ . Assumption 2 is satisfied if  $q_f(\bar{\gamma}) > \sqrt{\sigma/\beta}$ . This demand satisfies condition (5) with  $a_0 = 1$  and  $b_0 = -\mu$ . Condition (7) is satisfied in this example if

$$\frac{f'(\gamma)}{f(\gamma)} \geq \frac{2(1 - \alpha)}{\mu - \bar{\gamma}}$$

for all  $\gamma \in [\underline{\gamma}, \bar{\gamma}]$ .

**Constant elasticity demand.** Consider  $P(q) = q^{-\frac{1}{\eta}}$  with  $\eta > 1$ , and let  $Q = [0, q_{max}]$  where  $q_{max} > 0$ . For this example,  $q_f(\gamma) = \left(\frac{\gamma\eta}{\eta-1}\right)^{-\eta}$ ,  $v(q) = \frac{1}{\eta-1} q^{\frac{\eta-1}{\eta}}$  and  $\kappa = 1 + \frac{1}{\alpha} \frac{1}{\eta-1}$ . Assumption 1 is satisfied if  $q_{max}^{-\frac{1}{\eta}} < \frac{\gamma}{1-\frac{1}{\eta}(1-\frac{1}{\alpha})}$  where  $0 < \frac{\gamma}{1-\frac{1}{\eta}(1-\frac{1}{\alpha})} \leq \underline{\gamma}$  follows from  $\alpha \in (0, 1]$  and  $\underline{\gamma} > 0$ . Assumption 2 is satisfied if  $\left(\frac{\bar{\gamma}\eta}{\eta-1}\right)^{1-\eta} \frac{1}{\eta} > \sigma$ . This demand satisfies condition (5) with  $a_0 = -\frac{1}{\eta}$  and  $b_0 = 0$ . Condition (7) is satisfied in this example if

$$\frac{f'(\gamma)}{f(\gamma)} \geq -\frac{\alpha(\eta - 1) + 2}{\gamma}$$

for all  $\gamma \in [\underline{\gamma}, \bar{\gamma}]$ .

**Logarithmic demand.** Consider  $P(q) = \mu - \beta \ln q$  with  $\beta > 0$  and  $Q = [0, e^{\mu/\beta} - \epsilon]$  for  $\epsilon > 0$  small. For this example,  $q_f(\gamma) = e^{\frac{\mu-\beta-\gamma}{\beta}}$ ,  $v(q) = \beta q$  and  $\kappa = 1$ . Assumption 1 is satisfied for  $\epsilon > 0$  sufficiently small if  $\beta(1 - \alpha)/\alpha < \underline{\gamma}$ . Assumption 2 is satisfied if  $\beta e^{\frac{\mu-\beta-\bar{\gamma}}{\beta}} > \sigma$ . This demand satisfies condition (5) with  $a_0 = 0$  and  $b_0 = -\beta$ . Condition (7) is satisfied in this example if

$$\frac{f'(\gamma)}{f(\gamma)} \geq -\frac{\alpha}{\beta}$$

for all  $\gamma \in [\underline{\gamma}, \bar{\gamma}]$ .

Of course, the demand family defined by (5) includes examples beyond the three examples highlighted here.<sup>24</sup> The three examples, however, are commonly used in the literature and illustrate the breadth of the demand family defined by (5).

The sufficient conditions derived for the three examples admit an interpretation that is in line with the intuition developed previously whereby a rising density  $f(\gamma)$  and a concave  $v(q)$  work in favor of the optimality of the cap allocation. For the constant elasticity and log demand examples,  $v(q)$  is concave and linear, respectively, and the sufficient conditions hold when  $f(\gamma)$  is non-decreasing; indeed, for these examples, the sufficient conditions are satisfied even when  $f(\gamma)$  is decreasing, provided that it does not fall too quickly. By contrast, for the linear demand example,  $v(q)$  is convex, which works against the optimality of the cap allocation. The sufficient condition for this example thus places a more demanding restriction on the density: the condition fails if  $f(\gamma)$  is anywhere decreasing, and it requires that  $f(\gamma)$  is increasing (non-decreasing) when  $\alpha < 1$  ( $\alpha = 1$ ).

Interestingly, the demand family we have identified corresponds to the “linear delegation” case studied in [Kolotilin and Zapechelnnyuk \(2019\)](#) for the regulation problem when  $\sigma = 0$ .<sup>25</sup> However, we are not restricted to demand functions within the family that satisfies condition (5). For other demand functions, we could use parts (i) and (ii) of Propositions 1 and 2 directly. Alternatively, Corollary 1 also allows us to find simple conditions. Consider the following example, which does not fit in the family specified by (5):

**Exponential demand.** Consider  $P(q) = \beta e^{-q}$  with  $\beta > \max\{\bar{\gamma}, \underline{\gamma}e^2\}$  with  $Q = [0, 2 - \epsilon]$  for  $\epsilon > 0$  small. For this example, the sign of  $v''(q)$  varies over  $Q$ . We find that  $v''(q) = \beta e^{-q}(1 - q)$  and  $\kappa = 1 - \frac{1}{2\alpha}$ . Assumption 1 is satisfied for  $\epsilon > 0$  sufficiently small if  $\alpha > 2/(1 + \underline{\gamma})$  where this inequality when combined with  $\alpha \in (0, 1]$  implies that  $\underline{\gamma} > 1$ .

<sup>24</sup>For example, the demand function  $P(q) = \mu - \beta q^\eta$  satisfies (5) with  $a_0 = \eta$  and  $b_0 = -\mu\eta$ .

<sup>25</sup>[Kolotilin and Zapechelnnyuk \(2019\)](#) consider “linear delegation” problems where the principal’s objective,  $V(\gamma, q)$ , satisfies  $V_q(\gamma, q) = -\gamma - c(q)$  and where the agent’s objective,  $U(\gamma, q)$ , satisfies  $U_q(\gamma, q) = d(\gamma) - c(q)$  where  $c$  and  $d$  are continuous functions and  $c$  is strictly increasing. This implies that  $V_q(\gamma, q) - U_q(\gamma, q) = \gamma - d(\gamma)$ . That is,  $V_q(\gamma, q) - U_q(\gamma, q)$  is independent of  $q$ . In our case, for  $\sigma = 0$ ,  $V_q(\gamma, q) = w_q(\gamma, q) = -\gamma + b'(q) + \frac{1}{\alpha}v'(q)$ . Given that the objective of the agent can be modified by any strictly increasing affine transformation, we have that in our case,  $U_q(\gamma, q) = A(-\gamma + b'(q))$  for any  $A > 0$ . Hence, “linear delegation” requires that there exists  $A > 0$  such that  $V_q(\gamma, q) - U_q(\gamma, q)$  is independent of  $q$ , or alternatively, that there exists  $A > 0$  and  $B$  such that

$$b'(q) + \frac{1}{\alpha}v'(q) - Ab'(q) = B.$$

Using that  $b'(q) = P'(q)q + P(q)$  and that  $v'(q) = -qP'(q)$ , the above delivers condition (5). Note also that demand functions within this family deliver payoff functions  $w(\gamma, q)$  and  $b(q)$  that belong to the restricted preference family previously identified by [Amador and Bagwell \(2013\)](#) in their Proposition 2.

Assumption 2 is satisfied when  $\max_{q \in Q}(\beta e^{-q} - \bar{\gamma})q > \sigma$ . Corollary 1 holds if  $\alpha = 1$  and  $f$  is non-decreasing for all  $\gamma \in [\underline{\gamma}, \bar{\gamma}]$ .

At this point, it is convenient to pause and consider the “no-exclusion” scenario mentioned in the Introduction, wherein the regulator must ensure that all types choose to produce so that exclusion never occurs. This scenario can be motivated with reference to market settings where the monopolist provides an essential service with poor substitution alternatives. To characterize the optimal regulatory policy for the no-exclusion scenario, we may refer to the truncated cap allocation  $q_t^*(\gamma|\gamma_t)$ , defined in (3), for the special case where  $\gamma_t = \bar{\gamma}$ . This allocation corresponds to a price-cap regulatory policy, where the price cap is set at the second-best level that leaves a monopolist with the highest possible cost,  $\bar{\gamma}$ , with zero profit (inclusive of the fixed cost,  $\sigma$ ). To establish conditions for the optimality of this policy for the no-exclusion scenario, we simply set  $\gamma_t = \bar{\gamma}$  and refer to Proposition 1, Corollaries 1 and 2, and the demand examples above. Thus, for example, this price-cap allocation is optimal for the no-exclusion scenario if the demand function takes a linear, constant elasticity or log form and if a simple inequality condition holds, respectively, where the inequality condition is sure to hold if the density is non-decreasing over the full support.

By contrast, a characterization of optimal regulation for the general scenario in which exclusion is allowed must also determine the optimal value for  $\gamma_t$ . In other words, the optimal regulatory policy for the general scenario must determine as well the degree (if any) of exclusion. As mentioned in the Introduction, this scenario can be motivated with reference to market settings in which the monopolist provides an inessential service for a given market or region. We develop our results for the optimal degree of exclusion in the next section.

## 6 When to Exclude?

In the previous section, we obtain conditions that guarantee that the cap allocation with exclusion,  $q^*$  defined in (2), is optimal within the set of all feasible allocations. In this section, we study the properties of this optimal cap allocation,  $q^*$ , and in particular, whether or not some types are excluded from production.

Given a level of exclusion, parameterized by  $\gamma_t$ , we can write the welfare function as:

$$W(\gamma_t) = \int_{\underline{\gamma}}^{\gamma_H(q_i(\gamma_t))} (w(\gamma, q_f(\gamma)) - \sigma) dF(\gamma) + \int_{\gamma_H(q_i(\gamma_t))}^{\gamma_t} (w(\gamma, q_i(\gamma_t)) - \sigma) dF(\gamma)$$

where as before  $q_i(\gamma_t)$  represents the quantity strictly above  $q_f(\gamma_t)$  that makes type  $\gamma_t$  indifferent between producing or not.



Taking the derivative of the welfare function with respect to  $\gamma_t$ , we obtain that<sup>26</sup>

$$W'(\gamma_t) = (w(\gamma_t, q_i(\gamma_t)) - \sigma)f(\gamma_t) + \int_{\gamma_H(q_i(\gamma_t))}^{\gamma_t} w_q(\gamma, q_i(\gamma_t))q'_i(\gamma)dF(\gamma)$$

Given that  $q_i(\gamma_t)$  satisfies

$$\gamma_t = P(q_i(\gamma_t)) - \frac{\sigma}{q_i(\gamma_t)},$$

it follows that

$$q'_i(\gamma_t) = \frac{1}{P'(q_i(\gamma_t)) + \sigma/(q_i(\gamma_t))^2}$$

We have the following result:

**Lemma 5.** *The quantity of the indifferent type,  $q_i(\gamma_t)$ , is such that  $q'_i(\gamma_t) < 0$ . In addition,  $\gamma_t > \gamma_H(q_i(\gamma_t))$  for all  $\gamma_t > \underline{\gamma}$ .*

*Proof.* Given that  $q_i(\gamma_t) > q_f(\gamma_t)$ , it follows that  $\pi_q(\gamma_t, q_i(\gamma_t)) < 0$ . Hence,  $\pi_q(\gamma_t, q_i(\gamma_t)) = P'(q_i(\gamma_t))q_i(\gamma_t) + \pi(\gamma_t, q_i(\gamma_t))/q_i(\gamma_t) < 0$ . Using that  $\pi(\gamma_t, q_i(\gamma_t)) = \sigma$ , we obtain the first result of the lemma.

For the second result, there are two cases to consider, one where  $\gamma_H(\gamma_t) > \underline{\gamma}$  and the other where  $\gamma_H(\gamma_t) = \underline{\gamma}$ . For the latter case, the result is immediate. For the former case, we have that  $\gamma_H(\gamma_t) = b'(q_f(\gamma_H(\gamma_t))) = b'(q_i(\gamma_t)) = P'(q_i(\gamma_t))q_i(\gamma_t) + P(q_i(\gamma_t))$ . Thus,

$$\gamma_t - \gamma_H(\gamma_t) = - \left( \frac{\sigma}{q_i(\gamma_t)} + P'(q_i(\gamma_t))q_i(\gamma_t) \right)$$

The first result of the lemma establishes that the bracketed expression is negative; thus, it follows that  $\gamma_t > \gamma_H(\gamma_t)$ .  $\square$

We can use the definition of  $w$ , together with the definitions of  $\gamma_t$  and  $\gamma_H(\gamma_t)$ , to obtain that

$$\begin{aligned} W'(\gamma_t) &= \frac{1}{\alpha}v(q_i(\gamma_t))f(\gamma_t) - q_i(\gamma_t) \left( \frac{1}{\alpha}v'(q_i(\gamma_t)) \right) \frac{F(\gamma_t) - F(\gamma_H(\gamma_t))}{\gamma_t - b'(q_i(\gamma_t))} \\ &\quad + \frac{q_i(\gamma_t)}{\gamma_t - b'(q_i(\gamma_t))} \int_{\gamma_H(\gamma_t)}^{\gamma_t} (\gamma - b'(q_i(\gamma_t)))dF(\gamma) \end{aligned}$$

for all  $\gamma_t \in (\underline{\gamma}, \bar{\gamma}]$ .

---

<sup>26</sup> The function  $\gamma_H(\gamma_t)$  may fail to be differentiable at the highest value for  $\gamma_t$  at which  $\gamma_H(\gamma_t) = \underline{\gamma}$ ; however, the differentiability of  $\gamma_H$  does not affect the differentiability of the objective. An argument similar to the one use in the proof of Lemma 2 can be used to show differentiability of the objective.

We know that  $\gamma > b'(q_i(\gamma_t))$  for all  $\gamma > \gamma_H(\gamma_t)$ . So the last term in the above is strictly positive. Thus,

$$\begin{aligned} W'(\gamma_t) &> \frac{1}{\alpha} v(q_i(\gamma_t)) f(\gamma_t) - q_i(\gamma_t) \left( \frac{1}{\alpha} v'(q_i(\gamma_t)) \right) \frac{F(\gamma_t) - F(\gamma_H(\gamma_t))}{\gamma_t - b'(q_i(\gamma_t))} \\ &= \frac{1}{\alpha} v(q_i(\gamma_t)) \left[ f(\gamma_t) - \frac{F(\gamma_t) - F(\gamma_H(\gamma_t))}{\gamma_t - b'(q_i(\gamma_t))} \right] \\ &\quad + \frac{1}{\alpha} \left[ v(q_i(\gamma_t)) - v'(q_i(\gamma_t)) q_i(\gamma_t) \right] \frac{F(\gamma_t) - F(\gamma_H(\gamma_t))}{\gamma_t - b'(q_i(\gamma_t))}. \end{aligned}$$

If  $f(\gamma)$  is non-decreasing for all  $\gamma$ , we have that  $f(\gamma_t) - \frac{F(\gamma_t) - F(\gamma_H(\gamma_t))}{\gamma_t - \gamma_H(\gamma_t)} = \frac{\int_{\gamma_H(\gamma_t)}^{\gamma_t} [f(\gamma_t) - f(\gamma)] d\gamma}{\gamma_t - \gamma_H(\gamma_t)} \geq 0$ . Given that  $b'(q_i(\gamma_t)) \leq \gamma_H(\gamma_t) < \gamma_t$ , it follows that  $f(\gamma_t) - \frac{F(\gamma_t) - F(\gamma_H(\gamma_t))}{\gamma_t - b'(q_i(\gamma_t))} \geq 0$ . If  $v(q)$  is weakly concave, then  $v(q) - v'(q)q \geq 0$  as  $v(0) = 0$ . Hence:

**Proposition 3** (No exclusion). *If  $f(\gamma)$  is non-decreasing for all  $\gamma \in [\underline{\gamma}, \bar{\gamma}]$  and  $v(q)$  is weakly concave for all  $q \in Q$ , then  $W'(\gamma_t) > 0$  for all  $\gamma_t \in (\underline{\gamma}, \bar{\gamma}]$ . If the assumption of Proposition 2 holds, then,  $\gamma_t = \bar{\gamma}$  is optimal.*

*Proof.* The proof is given in the text. □

Proposition 3 delivers a general set of conditions under which there is no exclusion. The log demand and constant elasticity demand examples satisfy the requirement that  $v$  is weakly concave. In addition, if  $f$  is non-decreasing, then we may refer to Corollary 2 to conclude that the assumption of Proposition 2 holds in both examples; thus, the optimal allocation is the cap allocation without exclusion.

For the linear demand example, however,  $v$  is strictly convex, and Proposition 3 thus does not apply. For this example, in the case of a uniform distribution with  $\alpha = 1$ , we have a very different result:

**Proposition 4** (Exclusion). *Consider the linear demand example, suppose that  $F$  is uniform and  $\alpha = 1$ . If  $\sigma > 0$ , then*

- (a) *In any optimal allocation,  $\gamma_t$  is such that  $\gamma_H(\gamma_t) = \underline{\gamma}$ .*
- (b) *If  $q_i(\bar{\gamma}) < q_f(\underline{\gamma})$ , then in any optimal allocation  $\gamma_t < \bar{\gamma}$  and  $q_i(\gamma_t) = q^*$  where  $q^*$  is a solution of*

$$\frac{(\mu - \underline{\gamma})(\mu - \underline{\gamma} - 2\beta q^*)(q^*)^2}{\sigma - \beta(q^*)^2} = \sigma.$$

*Proof.* First, note that  $q'_i(\gamma_t) < 0$  implies that  $\gamma_H(\gamma_t)$  is strictly increasing in  $\gamma_t$ , as long as  $\gamma_H(\gamma_t) > \underline{\gamma}$ . That is, there exists a  $\hat{\gamma} \in (\underline{\gamma}, \bar{\gamma}]$  such that  $\gamma_H(\gamma_t) = \underline{\gamma}$  for all  $\gamma_t \leq \hat{\gamma}$  and

$\gamma_H(\gamma_t) > \underline{\gamma}$  for all  $\gamma_t > \hat{\gamma}$ . It is possible that  $\hat{\gamma} = \bar{\gamma}$ , and thus for any level of exclusion, all types are pooled.

Consider a situation where  $\hat{\gamma} < \bar{\gamma}$ . Then, for  $\gamma_t \in (\hat{\gamma}, \bar{\gamma}]$ , using the functional forms,  $\alpha = 1$ , the uniform distribution assumption, and that  $b'(q_i(\gamma_t)) = \gamma_H(\gamma_t)$ , we have that

$$\begin{aligned} W'(\gamma_t)/f_U &= v(q_i(\gamma_t)) - v'(q_i(\gamma_t))q_i(\gamma_t) + \frac{q_i(\gamma_t)}{\gamma_t - \gamma_H(\gamma_t)} \int_{\gamma_H(\gamma_t)}^{\gamma_t} (\gamma - \gamma_H(\gamma_t))d\gamma \\ &= -\beta q_i(\gamma_t)^2/2 - q_i(\gamma_t)\gamma_H(\gamma_t) + \frac{q_i(\gamma_t)}{2(\gamma_t - \gamma_H(\gamma_t))} (\gamma_t^2 - \gamma_H(\gamma_t)^2) \\ &= \frac{q_i(\gamma_t)}{2} [-\beta q_i(\gamma_t) + \gamma_t - \gamma_H(\gamma_t)] \end{aligned}$$

where  $f_U = \frac{1}{\bar{\gamma} - \underline{\gamma}}$  denotes the uniform density.

Using that  $\gamma_t - \gamma_H(\gamma_t) = -\frac{\sigma}{q_i(\gamma_t)} + \beta q_i(\gamma_t)$ , we obtain that

$$W'(\gamma_t)/f_U = -\frac{1}{2}\sigma < 0.$$

Thus, for all  $\gamma_t \in (\hat{\gamma}, \bar{\gamma}]$ ,  $W'(\gamma_t) < 0$ , and thus, in an optimal allocation  $\gamma_t \leq \hat{\gamma}$ , guaranteeing that all types are pooled. This completes the proof of part (a).

For part (b), note that the conditions imply that  $\hat{\gamma} < \bar{\gamma}$ . To see this, note that if  $\gamma_t = \bar{\gamma}$ , then  $q_i(\gamma_t) = q_i(\bar{\gamma}) < q_f(\underline{\gamma})$  under the conditions in part (b). There then exists  $\gamma_0 > \underline{\gamma}$  such that  $q_f(\gamma_0) = q_i(\gamma_t)$ , from which it follows that  $\gamma_H(\gamma_t) = \gamma_0 > \underline{\gamma}$ , a contradiction to part (a). Given that  $\gamma_t \leq \hat{\gamma} < \bar{\gamma}$ , part (a) implies that some types will be excluded.

For the case where  $\gamma_t < \hat{\gamma}$ , then  $\gamma_H(\gamma_t) = \underline{\gamma}$ , and

$$W'(\gamma_t)/f_U = -\frac{1}{2} \left( \sigma - \frac{(\mu - \underline{\gamma})(\mu - \underline{\gamma} - 2\beta q_i(\gamma_t))(q_i(\gamma_t))^2}{\sigma - \beta(q_i(\gamma_t))^2} \right).$$

Note that

$$\lim_{\gamma_t \rightarrow \underline{\gamma}} W'(\gamma_t)/f_U = -\frac{1}{2} \left( \sigma - \frac{(\mu - \underline{\gamma})(\mu - \underline{\gamma} - 2\beta q_i(\underline{\gamma}))(q_i(\underline{\gamma}))^2}{\sigma - \beta(q_i(\underline{\gamma}))^2} \right).$$

By definition of  $q_i$ , we have that  $(\mu - \beta q_i(\underline{\gamma}) - \underline{\gamma})q_i(\underline{\gamma}) = \sigma$ . Using this, we get

$$(\mu - \underline{\gamma} - 2\beta q_i(\underline{\gamma}))q_i(\underline{\gamma}) = \sigma - \beta(q_i(\underline{\gamma}))^2$$

and thus

$$\begin{aligned}\lim_{\gamma_t \rightarrow \underline{\gamma}} W'(\gamma_t)/f_U &= -\frac{1}{2} \left( \sigma - \frac{(\mu - \underline{\gamma})(\mu - \underline{\gamma} - 2\beta q_i(\underline{\gamma}))(q_i(\underline{\gamma}))^2}{\sigma - \beta(q_i(\underline{\gamma}))^2} \right) = -\frac{1}{2} (\sigma - (\mu - \underline{\gamma})q_i(\underline{\gamma})) \\ &= -\frac{1}{2} ((\mu - \beta q_i(\underline{\gamma}) - \underline{\gamma})q_i(\underline{\gamma}) - (\mu - \underline{\gamma})q_i(\underline{\gamma})) = \frac{1}{2} \beta q_i(\underline{\gamma})^2 > 0.\end{aligned}$$

We have already shown that  $W'(\gamma) < 0$  for  $\gamma \in (\hat{\gamma}, \bar{\gamma}]$ . Together with the above limiting result, and that  $W'$  is continuous at  $\hat{\gamma}$  (see footnote 26), it follows that the optimal  $\gamma_t$  is interior in  $[\underline{\gamma}, \hat{\gamma}]$  and thus satisfies  $W'(\gamma_t) = 0$ .  $\square$

This proposition contains two results. The first is that in the linear demand example with a uniform distribution and  $\alpha = 1$ , it is always optimal to pool all types at the cap (part (a)). Part (b) argues that if not all types pool at the cap when an allocation features no exclusion, that is, when  $\gamma_t = \bar{\gamma}$ , then some higher-cost types will necessary be excluded in any optimal allocation.<sup>27</sup>

The above result demonstrates that no-exclusion result of Proposition 3 is not a general property. Because of its tractability, the linear demand example with a uniform distribution and  $\alpha = 1$  is often used in the literature. For this case, we have shown that a cap allocation is optimal but that such an allocation also features the exclusion of higher-cost types.

## 7 Conclusion

We analyze the [Baron and Myerson \(1982\)](#) model of regulation under the restriction that transfers are infeasible. To do this, we extend the Lagrangian approach to delegation problems of [Amador and Bagwell \(2013\)](#) to include an ex post participation constraint that allows for the possible exclusion of some types. We report sufficient conditions under which optimal regulation takes the simple and common form of price-cap regulation. We identify families of demand and distribution functions and welfare weights that satisfy our sufficient conditions. We also report conditions under which the optimal price cap is set at a level such that no types are excluded. Using a linear demand example, we show that exclusion of higher-cost types can be optimal when these conditions fail to hold. Our analysis also can be used to provide conditions for the optimality of price-cap regulation when an ex post participation constraint is present and exclusion is infeasible.

---

<sup>27</sup> The proposition only characterizes the solution for  $\sigma > 0$ . When  $\sigma = 0$ , if  $q_i(\bar{\gamma}) < q_f(\underline{\gamma})$ , we can show that any  $\gamma_t$  such that  $q_i(\gamma_t) \leq q_f(\underline{\gamma})$  is optimal. Thus, the regulator is indifferent between some exclusion or none.

Our analysis points to several directions for future research. We mention four possibilities here.

First, we have provided general sufficient conditions so that a cap allocation with possible exclusion is optimal. These sufficient conditions guarantee that the Lagrangian approach can be used to show that a price cap is optimal *for any given level of exclusion*. Thus, the sufficient conditions may be stronger than necessary since the price-cap structure is required to be optimal even for exclusion levels that are sub-optimal. It should be possible to relax these conditions by using the Lagrangian approach *only* at the optimal level of exclusion.<sup>28</sup>

Second, when our sufficient conditions fail, it may be that the optimal allocation is not a price cap with possible exclusion. In that case, the Lagrangian approach requires us to identify the alternative solution candidate. It should be possible as well to construct Lagrange multipliers and generate sufficient conditions for such a case.

Third, we focus on a single-product monopolist and leave for future research the multi-product expression of our findings. More generally, the characterization of optimal delegation contracts in multi-dimensional settings is a challenging and important avenue for future work.<sup>29</sup>

Finally, our analysis extends the optimal delegation literature to include an ex-post participation constraint that allows for possible exclusion within a regulation framework. Many other applications may arise naturally in other environments, and they may be naturally captured in versions of the delegation model we have developed here.

---

<sup>28</sup>At the same time, a valuable by-product of our approach is that we obtain conditions under which a price cap is optimal for the no-exclusion scenario.

<sup>29</sup>For related work, see [Ambrus and Egorov \(2017\)](#), [Amador and Bagwell \(2018\)](#), [Armstrong and Vickers \(2010\)](#), [Frankel \(2014\)](#), [Frankel \(2016\)](#) and [Koessler and Martimort \(2012\)](#). The paper by [Frankel \(2016\)](#) is perhaps of special relevance here. He considers a model with multiple actions and establishes the exact optimality of a generalized cap rule, but under the assumptions that the loss function is quadratic, the agent has a constant bias, the ex ante distribution of states is normal iid, and the participation constraint is absent.

## A Proof of Lemma 2

*Proof.* We start by observing that for any  $\Delta \neq 0$ ,

$$\begin{aligned} & \frac{W^c(x + \Delta) - W^c(x)}{\Delta} \\ &= \int_{\gamma_c(x+\Delta)}^{\gamma_c(x)} \frac{w(\gamma, x + \Delta) - w(\gamma, q_f(\gamma))}{\Delta} dF(\gamma) + \int_{\gamma_c(x)}^{\bar{\gamma}} \frac{w(\gamma, x + \Delta) - w(\gamma, x)}{\Delta} dF(\gamma) \quad (8) \end{aligned}$$

Then, we consider two different cases.

**Case 1.**  $x > q_f(\underline{\gamma})$ . Then for all  $|\Delta| > 0$  small enough, we have that  $x + \Delta > q_f(\underline{\gamma})$ , and as a result  $\gamma_c(x + \Delta) = \gamma_c(x) = \underline{\gamma}$ , and thus

$$\frac{W^c(x + \Delta) - W^c(x)}{\Delta} = \int_{\gamma_c(x)}^{\bar{\gamma}} \frac{w(\gamma, x + \Delta) - w(\gamma, x)}{\Delta} dF(\gamma)$$

Taking the limit as  $\Delta \rightarrow 0$ , we obtain

$$\frac{dW^c(x)}{dx} = \int_{\gamma_c(x)}^{\bar{\gamma}} w_q(\gamma, x) dF(\gamma)$$

**Case 2.**  $0 < q_f(\bar{\gamma}) < x \leq q_f(\underline{\gamma})$ . Consider a neighborhood  $U_x$  around  $x$  such that that  $0 \notin cl(U_x)$ . Let  $K_x = \max_{y \in cl(U_x)} |b'(y) + v'(y)/\alpha|$ . Assumption 1 guarantees that such  $K_x$  exists and is finite. The mean value theorem guarantees that  $\left| \frac{(b(y)+v(y)/\alpha)-(b(x)-v(x)/\alpha)}{y-x} \right| \leq K_x$ .

Note that  $q_f(\gamma_c(x)) = x$ , and that for  $|\Delta| > 0$  small enough,  $q_f(\gamma) \in U_x$  for  $\gamma \in [\gamma_c(x + |\Delta|), \gamma_c(x - |\Delta|)]$ , given that  $q_f$  and  $\gamma_c$  are continuous. Then,

$$\begin{aligned} & \int_{\gamma_c(x+\Delta)}^{\gamma_c(x)} \left| \frac{w(\gamma, x + \Delta) - w(\gamma, q_f(\gamma))}{\Delta} \right| dF(\gamma) \\ &= \int_{\gamma_c(x+\Delta)}^{\gamma_c(x)} \left| \frac{-\gamma(x + \Delta - q_f(\gamma))}{\Delta} + \right. \\ & \quad \left. \frac{(b(x + \Delta) + v(x + \Delta)/\alpha) - (b(q_f(\gamma)) - v(q_f(\gamma))/\alpha)}{\Delta} \right| dF(\gamma) \\ &\leq \int_{\gamma_c(x+\Delta)}^{\gamma_c(x)} \gamma \left| \frac{x + \Delta - q_f(\gamma)}{\Delta} \right| dF(\gamma) + \\ & \quad \int_{\gamma_c(x+\Delta)}^{\gamma_c(x)} \left| \frac{(b(x + \Delta) + v(x + \Delta)/\alpha) - (b(q_f(\gamma)) - v(q_f(\gamma))/\alpha)}{x + \Delta - q_f(\gamma)} \right| \\ & \quad \times \left| \frac{x + \Delta - q_f(\gamma)}{\Delta} \right| dF(\gamma) \end{aligned}$$

$$\begin{aligned}
&\leq \int_{\gamma_c(x+\Delta)}^{\gamma_c(x)} \gamma \left| \frac{x+\Delta - q_f(\gamma)}{\Delta} \right| dF(\gamma) + \int_{\gamma_c(x+\Delta)}^{\gamma_c(x)} K_x \left| \frac{x+\Delta - q_f(\gamma)}{\Delta} \right| dF(\gamma) \\
&= (\gamma + K_x) \int_{\gamma_c(x+\Delta)}^{\gamma_c(x)} \left| \frac{x+\Delta - q_f(\gamma)}{\Delta} \right| dF(\gamma) \\
&\leq (\gamma + K_x) \int_{\gamma_c(x+\Delta)}^{\gamma_c(x)} \left| \frac{x+\Delta - x}{\Delta} \right| dF(\gamma) = (\gamma + K_x) \int_{\gamma_c(x+\Delta)}^{\gamma_c(x)} dF(\gamma)
\end{aligned}$$

The steps above are immediate except for the last inequality. For this we use that if  $x < q_f(\underline{\gamma})$ , then for all sufficiently small  $\Delta$ ,  $q_f(\gamma_c(x+\Delta)) = x+\Delta$ . If  $x = q_f(\underline{\gamma})$ , then for  $\Delta > 0$ , the integral range is empty (and thus, the integral equals zero). For  $\Delta < 0$ , we still have that  $q_f(\gamma_c(x+\Delta)) = x+\Delta$ .

Now note that the last integral above tends to zero as  $\Delta$  goes to zero. And thus, taking the limit of (8) as  $\Delta \rightarrow 0$ , we obtain that for  $x > q_f(\bar{\gamma})$ :

$$\frac{dW^c(x)}{dx} = \int_{\gamma_c(x)}^{\bar{\gamma}} w_q(\gamma, x) dF(\gamma)$$

Note that

$$\left. \frac{dW^c(s)}{dx} \right|_{x=q_{max}} = \int_{\gamma_c(q_{max})}^{\bar{\gamma}} w_q(\gamma, q_{max}) dF(\gamma) < \int_{\underline{\gamma}}^{\bar{\gamma}} w_q(\underline{\gamma}, q_{max}) dF(\gamma) = w_q(\underline{\gamma}, q_{max}) < 0$$

where we use that  $\gamma_c(q_{max}) = \underline{\gamma}$  as  $q_{max} > q_f(\underline{\gamma})$  and that  $w_q(\underline{\gamma}, q_{max}) > w_q(\gamma, q_{max})$  for  $\gamma > \underline{\gamma}$  to show the first inequality. For the last inequality we use Assumption 1.

Note that  $w_q(\bar{\gamma}, q_f(\bar{\gamma})) > 0$ , as  $v'(q) > 0$ . Consider  $x_0 > q_f(\bar{\gamma})$  such that  $w_q(\bar{\gamma}, x_0) > 0$ . Such an  $x_0$  exists by continuity of  $w_q$ . Note that for all  $q_0 \in (q_f(\bar{\gamma}), x_0]$ ,  $\gamma_c(q_0) < \bar{\gamma}$  and  $w_q(\bar{\gamma}, q_0) \geq w_q(\bar{\gamma}, x_0) > 0$ , by weak concavity of  $w$ . It follows that  $0 < w_q(\bar{\gamma}, q_0) \leq w_q(\gamma, q_0)$  for all  $\gamma \in [\underline{\gamma}, \bar{\gamma}]$ . Hence for all  $q_0 \in (q_f(\bar{\gamma}), x_0]$  we have that

$$\left. \frac{dW^c(s)}{dx} \right|_{x=q_0} = \int_{\gamma_c(q_0)}^{\bar{\gamma}} w_q(\gamma, q_0) dF(\gamma) > 0$$

It follows then that the optimal value of  $x$  is interior to  $(q_f(\bar{\gamma}), q_{max}]$  and must solve the first order condition in the lemma.  $\square$

## B Proof of Proposition 1

*Proof.* We proceed as follows. First, we re-state the Regulator's Truncated problem by expressing the incentive compatibility constraints in their standard form as an integral equation



and a monotonicity requirement:<sup>30</sup>

$$\begin{aligned}
& \max_{q_t: \Gamma_t(\gamma_t) \rightarrow Q} \int_{\Gamma_t(\gamma_t)} (w(\gamma, q_t(\gamma)) - \sigma) dF(\gamma) \quad \text{subject to:} \\
& -\gamma q_t(\gamma) + b(q_t(\gamma)) - \sigma - \int_{\gamma}^{\gamma_t} q_t(\tilde{\gamma}) d\tilde{\gamma} = \bar{U}, \text{ for all } \gamma \in \Gamma_t(\gamma_t) \\
& q_t(\gamma) \text{ non-increasing, for all } \gamma \in \Gamma_t(\gamma_t) \\
& 0 \leq -\gamma q_t(\gamma) + b(q_t(\gamma)) - \sigma, \text{ for all } \gamma \in \Gamma_t(\gamma_t)
\end{aligned}$$

where  $\bar{U} \equiv -\gamma_t q_t(\gamma_t) + b(q_t(\gamma_t)) - \sigma$  is the profit enjoyed by the monopolist with the highest possible cost type in  $\Gamma_t(\gamma_t)$ .

Next, we follow [Amador and Bagwell \(2013\)](#) and re-write the incentive constraints as a set of two inequalities and embed the monotonicity constraint in the choice set of  $q_t(\gamma)$ . With the choice set for  $q_t(\gamma)$  defined as  $\Phi \equiv \{q_t | q_t : \Gamma_t(\gamma_t) \rightarrow Q; \text{ and } q_t \text{ non-increasing}\}$ , the regulator's truncated problem may now be stated in final form as follows:

$$\max_{q_t \in \Phi} \int_{\Gamma_t(\gamma_t)} (w(\gamma, q_t(\gamma)) - \sigma) dF(\gamma) \quad \text{subject to:} \tag{P'_t}$$

$$\gamma q_t(\gamma) - b(q_t(\gamma)) + \sigma + \int_{\gamma}^{\gamma_t} q_t(\tilde{\gamma}) d\tilde{\gamma} + \bar{U} \leq 0, \text{ for all } \gamma \in \Gamma_t(\gamma_t) \tag{9}$$

$$-\gamma q_t(\gamma) + b(q_t(\gamma)) - \sigma - \int_{\gamma}^{\gamma_t} q_t(\tilde{\gamma}) d\tilde{\gamma} - \bar{U} \leq 0, \text{ for all } \gamma \in \Gamma_t(\gamma_t) \tag{10}$$

$$\gamma q_t(\gamma) - b(q_t(\gamma)) + \sigma \leq 0, \text{ for all } \gamma \in \Gamma_t(\gamma_t) \tag{11}$$

Let  $\Lambda_1(\gamma)$  and  $\Lambda_2(\gamma)$  denote the (cumulative) multiplier functions associated with the two inequalities that define the incentive compatibility constraints in the final form of the regulator's truncated problem. The multiplier functions  $\Lambda_1(\gamma)$  and  $\Lambda_2(\gamma)$  are restricted to be non-decreasing in  $\Gamma_t(\gamma_t)$ . Letting  $\Lambda(\gamma) \equiv \Lambda_1(\gamma) - \Lambda_2(\gamma)$ , we can write the Lagrangian of the regulator's truncated problem as stated in [P'\\_t](#) as follows:

$$\begin{aligned}
\mathcal{L} = & \int_{\Gamma_t} w(\gamma, q_t(\gamma)) dF(\gamma) - \int_{\Gamma_t} \left( \int_{\gamma}^{\gamma_t} q_t(\tilde{\gamma}) d\tilde{\gamma} + \bar{U} + \gamma q_t(\gamma) - b(q_t(\gamma)) + \sigma \right) d\Lambda(\gamma) \\
& + \int_{\Gamma_t} \left( -\gamma q_t(\gamma) + b(q_t(\gamma)) - \sigma \right) d\Psi(\gamma),
\end{aligned}$$

where without loss of generality we have removed the constant  $\sigma$  in the first integral and where to save notation we have removed the dependence of  $\Gamma_t$  on  $\gamma_t$ . Notice that  $\Psi(\gamma)$  is

---

<sup>30</sup>See, for example, [Milgrom and Segal \(2002\)](#).

the multiplier for the ex post participation constraints.  $\Psi(\gamma)$  is also restricted to be non-decreasing.

We propose the following multipliers:

$$\Lambda(\gamma) = \begin{cases} 0 & ; \gamma = \underline{\gamma} \\ w_q(\gamma, q_f(\gamma))f(\gamma) & ; \gamma \in (\underline{\gamma}, \gamma_H(\gamma_t)) \\ A + \kappa(F(\gamma_t) - F(\gamma)) & ; \gamma \in [\gamma_H(\gamma_t), \gamma_t] \end{cases}$$

and

$$\Psi(\gamma) = \begin{cases} 0 & ; \gamma \in [\underline{\gamma}, \gamma_t) \\ A & ; \gamma = \gamma_t \end{cases}$$

where

$$A = \frac{1}{\gamma_t - \gamma_H(\gamma_t)} \left[ \int_{\gamma_H(\gamma_t)}^{\gamma_t} w_q(\gamma, q_i(\gamma_t))f(\gamma)d\gamma + \kappa(\gamma_H(\gamma_t) - b'(q_i(\gamma_t)))F(\gamma_t) \right]. \quad (12)$$

Note that while defining  $\Lambda(\gamma)$ , we allow for the possibility that  $\gamma_H(\gamma_t) = \underline{\gamma}$ , and the intermediate case in the definition then does not apply. This is the case where there is full pooling of all types.

We show below that the hypothesis of Proposition 1 guarantees that  $R(\gamma) \equiv \kappa F(\gamma) + \Lambda(\gamma)$  is non-decreasing; thus, we may write  $\Lambda(\gamma)$  as the difference between two non-decreasing functions,  $\Lambda_1(\gamma) = R(\gamma)$  and  $\Lambda_2(\gamma) = \kappa F(\gamma)$ .<sup>31</sup> We also require that  $A \geq 0$  as  $\Phi$  must be non-decreasing. We establish this inequality below.

We note that the cap allocation  $q_t^*(\gamma|\gamma_t)$  together with the proposed multipliers satisfy complementary slackness. The incentive compatibility constraints bind under the cap allocation, and  $\Psi(\gamma)$  is constructed to be zero whenever the participation constraint holds with slack.

When these multipliers are used, the Lagrangian becomes

$$\begin{aligned} \mathcal{L} = \int_{\Gamma_t} w(\gamma, q_t(\gamma))dF(\gamma) - \int_{\Gamma_t} \left( \int_{\gamma}^{\gamma_t} q_t(\tilde{\gamma})d\tilde{\gamma} + \bar{U} + \gamma q_t(\gamma) - b(q_t(\gamma)) + \sigma \right) d\Lambda(\gamma) \\ + \left( -\gamma_t q_t(\gamma_t) + b(q_t(\gamma_t)) - \sigma \right) A \end{aligned}$$

Recalling the definition of  $\bar{U}$  and using  $\Lambda(\underline{\gamma}) = 0$  and  $\Lambda(\gamma_t) = A$ , we can then write the

---

<sup>31</sup>For our analysis, only the difference between  $\Lambda_1(\gamma)$  and  $\Lambda_2(\gamma)$  matters, and so we need only show that there exists two non-decreasing functions,  $\Lambda_1(\gamma)$  and  $\Lambda_2(\gamma)$ , whose difference delivers  $\Lambda(\gamma)$ .

Lagrangian as

$$\mathcal{L} = \int_{\Gamma_t} w(\gamma, q_t(\gamma)) dF(\gamma) - \int_{\Gamma_t} \left( \int_{\gamma}^{\gamma_t} q_t(\tilde{\gamma}) d\tilde{\gamma} + \gamma q_t(\gamma) - b(q_t(\gamma)) + \sigma \right) d\Lambda(\gamma)$$

Integrating the Lagrangian by parts we get<sup>32</sup>

$$\mathcal{L} = \int_{\Gamma_t} \left( w(\gamma, q_t(\gamma)) f(\gamma) - \Lambda(\gamma) q_t(\gamma) \right) d\gamma + \int_{\Gamma_t} \left( -\gamma q_t(\gamma) + b(q_t(\gamma)) - \sigma \right) d\Lambda(\gamma) \quad (13)$$

Let us now consider the concavity of the Lagrangian. Using (13), we may re-write the Lagrangian as

$$\begin{aligned} \mathcal{L} &= \int_{\Gamma_t} \left( w(\gamma, q_t(\gamma)) - \kappa(-\gamma q_t(\gamma) + b(q_t(\gamma)) - \sigma) \right) f(\gamma) d\gamma - \int_{\Gamma_t} \Lambda(\gamma) q_t(\gamma) d\gamma \\ &\quad + \int_{\Gamma_t} \left( -\gamma q_t(\gamma) + b(q_t(\gamma)) - \sigma \right) d(\kappa F(\gamma) + \Lambda(\gamma)) \end{aligned}$$

From the definition of  $\kappa$ ,  $w(\gamma, q_t(\gamma)) - \kappa b(q_t(\gamma))$  is concave in  $q_t(\gamma)$ . We may thus conclude that the Lagrangian is concave in  $q_t(\gamma)$  if

$$\kappa F(\gamma) + \Lambda(\gamma)$$

is non-decreasing for all  $\gamma \in [\underline{\gamma}, \gamma_t]$ . Using the constructed  $\Lambda(\gamma)$  and referring to part (ii) of Proposition 1, we see that  $\kappa F(\gamma) + \Lambda(\gamma)$  is non-decreasing for all  $\gamma \in [\underline{\gamma}, \gamma_t]$  if the jumps in  $\Lambda(\gamma)$  at  $\underline{\gamma}$  and  $\gamma_H(\gamma_t)$  are non-negative. We verify these jumps are indeed non-negative below.

We now show that the cap allocation  $q_t^*$  maximizes the Lagrangian. To this end, we use the sufficiency part of Lemma A.2 in Amador et al. (2006), which concerns the maximization of concave functionals on a convex cone. In our case, we need to extend the set  $Q$  to be  $[0, \infty)$ , making our choice set  $\Phi$  a convex cone. To do this, we follow Amador and Bagwell (2013) and extend  $b$  and  $w$  to the entire non-negative ray of the real line. We can then apply Lemma A.2 to the extended Lagrangian with the choice set  $\widehat{\Phi} \equiv \{q|q : \Gamma_t \rightarrow \mathbb{R}_+; \text{ and } q \text{ non-increasing}\}$ .

---

<sup>32</sup>Observe that  $h(\gamma) \equiv \int_{\gamma}^{\gamma_t} q_t(\tilde{\gamma}) d\tilde{\gamma}$  exists (as  $q_t$  is bounded and measurable by monotonicity) and is absolutely continuous. Observe as well that  $\Lambda(\gamma) \equiv \Lambda_1(\gamma) - \Lambda_2(\gamma)$  is a function of bounded variation, as it is the difference between two non-decreasing and bounded functions. We may thus conclude that  $\int_{\gamma}^{\gamma_t} h(\gamma) d\Lambda(\gamma)$  exists (it is the Riemman-Stieltjes integral), and integration by parts can be done as follows:  $\int_{\underline{\gamma}}^{\gamma_t} h(\gamma) d\Lambda(\gamma) = h(\bar{\gamma})\Lambda(\bar{\gamma}) - h(\underline{\gamma})\Lambda(\underline{\gamma}) - \int_{\underline{\gamma}}^{\gamma_t} \Lambda(\gamma) dh(\gamma)$ . Given that  $h(\gamma)$  is absolutely continuous, we can replace  $dh(\gamma)$  with  $-q_t(\gamma)d\gamma$ .

Following the arguments in [Amador and Bagwell \(2013\)](#), we can then establish that the cap allocation  $q_t^*$  maximizes the Lagrangian if the Lagrangian is concave and the following first order conditions hold:

$$\begin{aligned}\partial\mathcal{L}(q_t^*; q_t^*) &= 0 \\ \partial\mathcal{L}(q_t^*; x) &\leq 0 \text{ for all } x \in \Phi,\end{aligned}$$

where  $\partial\mathcal{L}(q_t^*; x)$  is the Gateaux differential of the Lagrangian in (13) in the direction  $x$ .<sup>33</sup> Importantly, the Lagrangian in (13) is evaluated using our constructed multiplier functions.

Taking the Gateaux differential of the Lagrangian in (13) in direction  $x \in \Phi$ , we get<sup>34</sup>

$$\begin{aligned}\partial\mathcal{L}(q_t^*; x) &= \int_{\Gamma_t} \left( w_q(\gamma, q_t^*(\gamma))f(\gamma) - \Lambda(\gamma) \right) x(\gamma) d\gamma \\ &\quad + \int_{\Gamma_t} \left( -\gamma + b'(q_t^*(\gamma)) \right) x(\gamma) d\Lambda(\gamma).\end{aligned}$$

Using  $b'(q_f(\gamma)) = \gamma$  and our knowledge of  $\Lambda$  and  $\Psi$ , we get that

$$\partial\mathcal{L}(q_t^*; x) = \int_{\gamma_H(\gamma_t)}^{\gamma_t} \left( w_q(\gamma, q_i(\gamma_t))f(\gamma) - A - \kappa(F(\gamma_t) - F(\gamma)) - \kappa(b'(q_i(\gamma_t)) - \gamma)f(\gamma) \right) x(\gamma) d\gamma$$

Hence, integrating by parts, we get

$$\begin{aligned}\partial\mathcal{L}(q_t^*; x) &= \left[ \int_{\gamma_H(\gamma_t)}^{\gamma_t} \left( w_q(\gamma, q_i(\gamma_t))f(\gamma) - A - \kappa(F(\gamma_t) - F(\gamma)) - \kappa(b'(q_i(\gamma_t)) - \gamma)f(\gamma) \right) d\gamma \right] x(\gamma_t) \\ &\quad - \int_{\gamma_H(\gamma_t)}^{\gamma_t} \left[ \int_{\gamma_H(\gamma_t)}^{\gamma} \left( w_q(\tilde{\gamma}, q_i(\gamma_t))f(\tilde{\gamma}) - A - \kappa(F(\gamma_t) - F(\tilde{\gamma})) - \kappa(b'(q_i(\gamma_t)) - \tilde{\gamma})f(\tilde{\gamma}) \right) d\tilde{\gamma} \right] dx(\gamma)\end{aligned}$$

Now, we use  $\int_b^a \{(F(c) - F(x)) + (d - x)f(x)\} dx = (a - b)(F(c) - F(a)) + (d - b)(F(a) - F(b))$  to get that

$$\partial\mathcal{L}(q_t^*; x) = \left[ \int_{\gamma_H(\gamma_t)}^{\gamma_t} w_q(\gamma, q_i(\gamma_t))f(\gamma) d\gamma - (\gamma_t - \gamma_H(\gamma_t))A \right]$$

---

<sup>33</sup>Given a function  $T : \Omega \rightarrow Y$ , where  $\Omega \subset X$  and  $X$  and  $Y$  are normed spaces, if for  $x \in \Omega$  and  $h \in X$  the limit

$$\lim_{\alpha \downarrow 0} \frac{1}{\alpha} [T(x + \alpha h) - T(x)]$$

exists, then it is called the Gateaux differential at  $x$  with direction  $h$  and is denoted by  $\partial T(x; h)$ .

<sup>34</sup>Existence of the Gateaux differential follows from Lemma A.1 in [Amador et al. \(2006\)](#). See [Amador and Bagwell \(2013\)](#) for further details concerning the application of this lemma.

$$\begin{aligned}
& -\kappa(b'(q_i(\gamma_t)) - \gamma_H(\gamma_t))(F(\gamma_t) - F(\gamma_H(\gamma_t))) \Big] x(\gamma_t) \\
& - \int_{\gamma_H(\gamma_t)}^{\gamma_t} \left[ \int_{\gamma_H(\gamma_t)}^{\gamma} w_q(\tilde{\gamma}, q_i(\gamma_t)) f(\tilde{\gamma}) d\tilde{\gamma} - (\gamma - \gamma_H(\gamma_t))A \right. \\
& \left. - \kappa((\gamma - \gamma_H(\gamma_t))(F(\gamma_t) - F(\gamma)) + (b'(q_i(\gamma_t)) - \gamma_H(\gamma_t))(F(\gamma) - F(\gamma_H(\gamma_t)))) \right] dx(\gamma)
\end{aligned}$$

Given that  $(b'(q_i(\gamma_t)) - \gamma_H(\gamma_t))F(\gamma_H(\gamma_t)) = 0$ , as  $\gamma_H(\gamma_t) < \gamma_t$  and  $b'(q_i(\gamma_t)) = \gamma_H(\gamma_t)$  if  $\gamma_H(\gamma_t) \in (\underline{\gamma}, \bar{\gamma})$ , the above becomes:

$$\begin{aligned}
\partial \mathcal{L}(q_t^*; x) = & \left[ \int_{\gamma_H(\gamma_t)}^{\gamma_t} w_q(\gamma, q_i(\gamma_t)) f(\gamma) d\gamma - (\gamma_t - \gamma_H(\gamma_t))A \right. \\
& \left. + \kappa(\gamma_H(\gamma_t) - b'(q_i(\gamma_t)))F(\gamma_t) \right] x(\gamma_t) \\
& - \int_{\gamma_H(\gamma_t)}^{\gamma_t} \left[ \int_{\gamma_H(\gamma_t)}^{\gamma} w_q(\tilde{\gamma}, q_i(\gamma_t)) f(\tilde{\gamma}) d\tilde{\gamma} - (\gamma - \gamma_H(\gamma_t))A \right. \\
& \left. - \kappa(\gamma - \gamma_H(\gamma_t))F(\gamma_t) + \kappa(\gamma - b'(q_i(\gamma_t)))F(\gamma) \right] dx(\gamma)
\end{aligned}$$

Using the definition of  $G$  in equation (4), we can rewrite the above as

$$\partial \mathcal{L}(q_t^*; x) = (G(\gamma_t|\gamma_t) - A)(\gamma_t - \gamma_H(\gamma_t))x(\gamma_t) - \int_{\gamma_H(\gamma_t)}^{\gamma_t} (G(\gamma|\gamma_t) - A)(\gamma - \gamma_H(\gamma_t))dx(\gamma).$$

Using (4) and (12), we also observe that

$$G(\gamma_t|\gamma_t) = A \tag{14}$$

and thus

$$\partial \mathcal{L}(q_t^*; x) = - \int_{\gamma_H(\gamma_t)}^{\gamma_t} (G(\gamma|\gamma_t) - A)(\gamma - \gamma_H(\gamma_t))dx(\gamma). \tag{15}$$

We are now ready to evaluate the first order conditions.

Note that it follows immediately that  $\partial \mathcal{L}(q_t^*; q_t^*) = 0$  as  $q_t^*$  is constant for  $\gamma \in [\gamma_H(\gamma_t), \gamma_t]$ .

If  $G(\gamma|\gamma_t) \leq A = G(\gamma_t|\gamma_t)$  for all  $\gamma \in [\gamma_H(\gamma_t), \gamma_t]$ , then for any non-increasing  $x \in \Phi$ , it follows that  $\partial \mathcal{L}(q_t^*; x) \leq 0$ , which is provided by part (i) of Proposition 1.

Recall also that we require  $A \geq 0$ , since  $\Phi$  must be non-decreasing. To see that this

inequality holds, note that

$$A = \kappa \left[ \frac{\gamma_H(\gamma_t) - b'(q_i(\gamma_t))}{\gamma_t - \gamma_H(\gamma_t)} \right] F(\gamma_t) + \frac{1}{\gamma_t - \gamma_H(\gamma_t)} \int_{\gamma_H(\gamma_t)}^{\gamma_t} w_q(\tilde{\gamma}, q_i(\gamma_t)) f(\tilde{\gamma}) d\tilde{\gamma}.$$

By the definition of  $\gamma_H$ , we have that  $q_i(\gamma_t) \geq q_f(\gamma_H(\gamma_t))$ . Note also that  $b'(q_f(\gamma_H(\gamma_t))) = \gamma_H(\gamma_t)$ , and concavity of  $b$  implies that  $b'(q_i(\gamma_t)) \leq b'(q_f(\gamma_H(\gamma_t))) = \gamma_H(\gamma_t)$ . So the first term in the previous equation is non-negative. Finally, note that  $w_q(\gamma, q) = P(q) - \gamma + P'(q)q + \frac{1}{\alpha}v'(q) = P(q) - \gamma - \frac{1-\alpha}{\alpha}qP'(q)$ . Thus

$$\begin{aligned} w_q(\gamma_t, q_i(\gamma_t)) &= P(q_i(\gamma_t)) - \gamma_t - \frac{1-\alpha}{\alpha}q_i(\gamma_t)P'(q_i(\gamma_t)) \\ &= \frac{\sigma}{q_i(\gamma_t)} - \frac{1-\alpha}{\alpha}q_i(\gamma_t)P'(q_i(\gamma_t)) \geq 0 \end{aligned}$$

where the last equality follows from  $(P(q_i(\gamma_t)) - \gamma_t)q_i(\gamma_t) = \sigma$ , by the definition of  $q_i$ . But we also have that

$$w_q(\gamma, q) > w_q(\gamma', q)$$

for all  $\gamma < \gamma'$ , and thus

$$w_q(\gamma, q_i(\gamma_t)) \geq w_q(\gamma_t, q_i(\gamma_t)) \geq 0$$

for  $\gamma \leq \gamma_t$ . Hence, we can also sign the integral term:  $\int_{\gamma_H(\gamma_t)}^{\gamma_t} w_q(\tilde{\gamma}, q_i(\gamma_t)) f(\tilde{\gamma}) d\tilde{\gamma} \geq 0$ . Taken together, the above implies that  $A \geq 0$ .

As discussed above, we now finish the argument that  $\kappa F(\gamma) + \Lambda(\gamma)$  is non-decreasing for all  $\gamma \in [\underline{\gamma}, \bar{\gamma}]$  by showing that the potential jumps in  $\Lambda(\gamma)$  are non-negative. There are two cases to consider. The first case is where  $\gamma_H(\gamma_t) > \underline{\gamma}$ . In this case, there are two jumps, one at  $\underline{\gamma}$  and one at  $\gamma_H(\gamma_t)$ . For the jump at  $\gamma_H(\gamma_t)$ , we get

$$A + \kappa(F(\gamma_t) - F(\gamma_H(\gamma_t))) - w_q(\gamma_H(\gamma_t), q_f(\gamma_H(\gamma_t)))f(\gamma_H(\gamma_t)) = G(\gamma_t|\gamma_t) - G(\gamma_H(\gamma_t)|\gamma_t)$$

where  $G(\gamma_H(\gamma_t)|\gamma_t) = -\kappa[F(\gamma_t) - F(\gamma_H(\gamma_t))] + w_q(\gamma_H(\gamma_t), q_f(\gamma_H(\gamma_t)))f(\gamma_H(\gamma_t))$ . Part (i) of Proposition 1 guarantees that  $G(\gamma_t|\gamma_t) \geq G(\gamma_H(\gamma_t)|\gamma_t)$ , and thus the jump at  $\gamma_H(\gamma_t)$  is non-negative.

The jump in  $\Lambda(\gamma)$  at  $\underline{\gamma}$  is non-negative, since  $w_q(\underline{\gamma}, q_f(\underline{\gamma}))f(\underline{\gamma}) > 0$ .

Finally, for the case where  $\gamma_H(\gamma_t) = \underline{\gamma}$ , there is only one jump, at  $\underline{\gamma}$ . The jump is

$$A + \kappa F(\gamma_t)$$

which is positive, given that we have shown that  $A \geq 0$ .

To complete the proof, we use Theorem 1 in [Amador and Bagwell \(2013\)](#). To apply this theorem, we set 1. (i)  $x_0 \equiv q_t^*$ ; (ii)  $X \equiv \{q_t | q_t : \Gamma_t \rightarrow Q\}$ ; (iii)  $f$  to be given by the negative of the objective function,  $\int_{\Gamma_t} w(\gamma, q_t(\gamma)) dF(\gamma)$ , as a function of  $q_t \in X$ ; (iv)  $Z \equiv \{(z_1, z_2, z_3) | z_1 : \Gamma_t \rightarrow \mathbb{R}, z_2 : \Gamma_t \rightarrow \mathbb{R} \text{ and } z_3 : \Gamma_t \rightarrow \mathbb{R} \text{ with } z_1, z_2, z_3 \text{ integrable}\}$ ; (v)  $\Omega \equiv \Phi$ ; (vi)  $P \equiv \{(z_1, z_2, z_3) | (z_1, z_2, z_3) \in Z \text{ such that } z_1(\gamma) \geq 0, z_2(\gamma) \geq 0 \text{ and } z_3(\gamma) \geq 0 \text{ for all } \gamma \in \Gamma_t\}$ ; (vii)  $\hat{G}$  (which is referred to as  $G$  in Theorem 1) to be the mapping from  $\Phi$  to  $Z$  given by the left hand sides of inequalities (9), (10) and (11); (viii)  $T$  to be the linear mapping:

$$T((z_1, z_2, z_3)) \equiv \int_{\Gamma_t} z_1(\gamma) d\Lambda_1(\gamma) + \int_{\Gamma_t} z_2(\gamma) d\Lambda_2(\gamma) + \int_{\Gamma_t} z_3(\gamma) d\Psi(\gamma)$$

where  $\Lambda_1$ ,  $\Lambda_2$  and  $\Psi$  being non-decreasing functions implies that  $T(z) \geq 0$  for  $z \in P$ . We have that

$$\begin{aligned} T(\hat{G}(x_0)) &\equiv \int_{\Gamma_t} \left( \int_{\gamma}^{\gamma_t} q_t^*(\tilde{\gamma}) d\tilde{\gamma} + \bar{U} + \gamma q_t^*(\gamma) - b(q_t^*(\gamma)) + \sigma \right) d(\Lambda_1(\gamma) - \Lambda_2(\gamma)) \\ &\quad - \int_{\Gamma_t} \left( -\gamma q_t^*(\gamma) + b(q_t^*(\gamma)) - \sigma \right) d\Psi(\gamma) = 0 \end{aligned}$$

where  $\bar{U}$  is evaluated at the  $q_t^*$  allocation, and where the last equality follows from the  $q_t^*$  allocation and the proposed multipliers. We have found conditions under which the proposed allocation,  $q_t^*$ , minimizes  $f(x) + T(\hat{G}(x))$  for  $x \in \Omega$ . Given that  $T(\hat{G}(x_0)) = 0$ , then the conditions of Theorem 1 hold and it follows that  $q_t^*$  solves  $\min_{x \in \Omega} f(x)$  subject to  $-\hat{G}(x) \in P$ , which is Problem  $P'_t$ .  $\square$

## B.1 Proof of Corollary 1

*Proof.* Letting  $q_i$  and  $\gamma_H$  represent  $q_i(\gamma_t)$  and  $\gamma_H(\gamma_t)$ , respectively, we start with the following manipulations:

$$\begin{aligned} G(\gamma | \gamma_t) &= -\kappa F(\gamma_t) + \kappa \frac{\gamma - b'(q_i)}{\gamma - \gamma_H} F(\gamma) + \frac{1}{\gamma - \gamma_H} \int_{\gamma_H}^{\gamma} \left( -\tilde{\gamma} + b'(q_i) + \frac{1}{\alpha} v'(q_i) \right) f(\tilde{\gamma}) d\tilde{\gamma} \\ &= -\kappa F(\gamma_t) + \kappa \frac{\gamma}{\gamma - \gamma_H} F(\gamma) - \kappa \frac{b'(q_i)}{\gamma - \gamma_H} F(\gamma) \\ &\quad + \kappa \frac{b'(q_i)}{\gamma - \gamma_H} F(\gamma_H) - \kappa \frac{b'(q_i)}{\gamma - \gamma_H} F(\gamma_H) + \frac{1}{\gamma - \gamma_H} \int_{\gamma_H}^{\gamma} \left( -\tilde{\gamma} + b'(q_i) + \frac{1}{\alpha} v'(q_i) \right) f(\tilde{\gamma}) d\tilde{\gamma} \\ &= -\kappa F(\gamma_t) + \frac{\kappa}{\gamma - \gamma_H} \int_{\gamma_H}^{\gamma} (\tilde{\gamma} f(\tilde{\gamma}) + F(\tilde{\gamma})) d\tilde{\gamma} - \kappa \frac{b'(q_i)}{\gamma - \gamma_H} (F(\gamma) - F(\gamma_H)) \end{aligned}$$

$$\begin{aligned}
& + \frac{1}{\gamma - \gamma_H} \int_{\gamma_H}^{\gamma} \left( -\tilde{\gamma} + b'(q_i) + \frac{1}{\alpha} v'(q_i) \right) f(\tilde{\gamma}) d\tilde{\gamma} \\
& = -\kappa F(\gamma_t) + \frac{\kappa}{\gamma - \gamma_H} \int_{\gamma_H}^{\gamma} (\tilde{\gamma} f(\tilde{\gamma}) + F(\tilde{\gamma}) - b'(q_i) f(\tilde{\gamma})) d\tilde{\gamma} \\
& \quad + \frac{1}{\gamma - \gamma_H} \int_{\gamma_H}^{\gamma} \left( -\tilde{\gamma} + b'(q_i) + \frac{1}{\alpha} v'(q_i) \right) f(\tilde{\gamma}) d\tilde{\gamma} \\
& = -\kappa F(\gamma_t) + \frac{1}{\gamma - \gamma_H} \int_{\gamma_H}^{\gamma} \left[ \kappa F(\tilde{\gamma}) + \frac{1}{\alpha} v'(q_i) f(\tilde{\gamma}) + (\kappa - 1)(\tilde{\gamma} - b'(q_i)) f(\tilde{\gamma}) \right] d\tilde{\gamma} \\
& = -\kappa F(\gamma_t) + \frac{1}{\gamma - \gamma_H} \int_{\gamma_H}^{\gamma} M_2(\tilde{\gamma}) d\tilde{\gamma},
\end{aligned}$$

where we use in the third equality above that  $\frac{\gamma_H - b'(q_i)}{\gamma - \gamma_H} F(\gamma_H) = 0$  and where we define

$$M_2(\tilde{\gamma}) \equiv \kappa F(\tilde{\gamma}) + \frac{1}{\alpha} v'(q_i) f(\tilde{\gamma}) + (\kappa - 1)(\tilde{\gamma} - b'(q_i)) f(\tilde{\gamma}).$$

Thus,

$$(\gamma - \gamma_H) G(\gamma | \gamma_t) = -\kappa(\gamma - \gamma_H) F(\gamma_t) + \int_{\gamma_H}^{\gamma} M_2(\tilde{\gamma}) d\tilde{\gamma}.$$

Taking a derivative with respect to  $\gamma$ , for  $\gamma > \gamma_H$ , we obtain

$$(\gamma - \gamma_H) G'(\gamma | \gamma_t) + G(\gamma | \gamma_t) = -\kappa F(\gamma_t) + M_2(\gamma)$$

and thus

$$(\gamma - \gamma_H) G'(\gamma | \gamma_t) = M_2(\gamma) - \frac{1}{\gamma - \gamma_H} \int_{\gamma_H}^{\gamma} M_2(\tilde{\gamma}) d\tilde{\gamma}.$$

It follows that, if  $M_2'(\gamma) \geq 0$ , then  $G'(\gamma | \gamma_t) \geq 0$ . Now note that

$$\begin{aligned}
M_2'(\gamma) & = \kappa f(\gamma) + \frac{1}{\alpha} v'(q_i) f'(\gamma) + (\kappa - 1)(\gamma - b'(q_i)) f'(\gamma) + (\kappa - 1) f(\gamma) \\
& = (2\kappa - 1) f(\gamma) + \kappa(\gamma - b'(q_i)) f'(\gamma) + (-\gamma + b'(q_i) + v'(q_i)/\alpha) f'(\gamma).
\end{aligned}$$

Recall that

$$\gamma - b'(q_i) \geq 0$$

for  $\gamma \geq \gamma_H$ . In addition,

$$-\gamma + b'(q_i) + v'(q_i)/\alpha = -\gamma + b'(q_i) + v'(q_i) + \left( \frac{1}{\alpha} - 1 \right) v'(q_i)$$



$$= (P(q_i) - \gamma) + \left(\frac{1}{\alpha} - 1\right) v'(q_i) \geq 0 \text{ for } \gamma \geq \gamma_H$$

where we use that  $b'(q_i) + v'(q_i) = P(q_i)$  and where the inequality follows from  $v'(q_i) > 0$ ,  $\alpha \in (0, 1]$ , and that  $P(q_i) \geq \gamma$  for all types in  $[\gamma_H, \gamma_t]$  (so that they can make profits and cover the fixed cost  $\sigma \geq 0$ ). Hence,

$$M'_2(\gamma) = (2\kappa - 1)f(\gamma) + \kappa(\text{non-negative term})f'(\gamma) + (\text{non-negative term})f'(\gamma).$$

Thus,  $\kappa \geq 1/2$  and  $f'(\gamma) \geq 0$  together are sufficient to guarantee that  $M'_2(\gamma) \geq 0$  and thus that  $G(\gamma|\gamma_t)$  is non-decreasing for any  $\gamma_t$ . Hence, part (i) of Proposition 1 then holds for all  $\gamma_t \in (\underline{\gamma}, \bar{\gamma}]$ .

Finally note that

$$M'_1(\gamma) = \kappa f(\gamma) + \frac{1}{\alpha} v''(q_f(\gamma)) q'_f(\gamma) f(\gamma) + \frac{1}{\alpha} v'(q_f(\gamma)) f'(\gamma).$$

Using  $q'_f(\gamma) = 1/b''(q_f(\gamma))$  and the definition of  $\kappa$ , we obtain that

$$M'_1(\gamma) \geq (2\kappa - 1)f(\gamma) + \frac{1}{\alpha} v'(q_f(\gamma)) f'(\gamma) \geq 0$$

where the second inequality follows from  $\kappa \geq 1/2$  and  $f$  non-decreasing. Thus part (ii) of Proposition 1 also holds for all  $\gamma_t \in (\gamma_t, \bar{\gamma}]$ . We can thus use Proposition 2 to obtain the desired result.  $\square$

## References

- Alonso, Ricardo and Niko Matouschek, “Optimal Delegation,” *The Review of Economic Studies*, 2008, 75 (1), 259–293. [7](#), [1](#)
- Amador, Manuel and Kyle Bagwell, “Tariff Revenue and Tariff Caps,” *American Economic Review, Papers and Proceedings*, May 2012, 102 (3), 459–465. [9](#)
- and —, “The theory of optimal delegation with an application to tariff caps,” *Econometrica*, 2013, 81 (4), 1541–1599. [\(document\)](#), [1](#), [11](#), [2](#), [13](#), [3.1](#), [4.1](#), [4.1](#), [25](#), [7](#), [B](#), [B](#), [34](#), [B](#)
- and —, “Money Burning in the Theory of Delegation,” Working paper, Stanford University 2018. [9](#), [11](#), [29](#)
- , Iván Werning, and George-Marios Angeletos, “Commitment vs. Flexibility,” *Econometrica*, 2006, 74 (2), 365–96. [1](#), [9](#), [B](#), [34](#)
- , Kyle Bagwell, and Alex Frankel, “A Note on Interval Delegation,” *Economic Theory Bulletin*, 2018, 6 (2), 239–49. [9](#)
- Ambrus, Attila and Georgy Egorov, “Delegation and Nonmonetary Incentives,” *Journal of Economic Theory*, 2017, 171, 101–135. [9](#), [11](#), [29](#)
- Armstrong, Mark, “Multiproduct Nonlinear Pricing,” *Econometrica*, 1996, 64 (1), 51–75. [6](#)
- and David E. M. Sappington, “Recent Developments in the Theory of Regulation,” in Mark Armstrong and Robert Porter, eds., *Handbook of Industrial Organization*, Vol. 3, North-Holland, Amsterdam, 2007, pp. 1557–1700. [1](#), [2](#), [5](#)
- and John Vickers, “A Model of Delegated Project Choice,” *Econometrica*, 1 2010, 78 (1), 213–244. [9](#), [29](#)
- Athey, Susan, Andrew Atkeson, and Patrick J. Kehoe, “The Optimal Degree of Discretion in Monetary Policy,” *Econometrica*, 2005, 73 (5), 1431–1475. [9](#)
- , Kyle Bagwell, and Chris Sanchirico, “Collusion and Price Rigidity,” *The Review of Economic Studies*, 2004, 71 (2), 317–349. [9](#)
- Baron, David P., “Design of Regulatory Mechanisms and Institutions,” in Richard Schmalensee and Robert D. Willig, eds., *Handbook of Industrial Organization*, Vol. 2, North-Holland, Amsterdam, 1989, pp. 1347–1447. [2](#)

- **and Roger B. Myerson**, “Regulating a Monopolist with Unknown Costs,” *Econometrica*, 1982, *50* (4), 911–930. [\(document\)](#), [1](#), [7](#), [7](#)
- Burkett, Justin**, “Optimally Constraining a Bidder Using a Simple Budget,” *Theoretical Economics*, 2016, *11*, 133–155. [9](#)
- Church, Jeffrey and Roger Ware**, *Industrial Organization: A Strategic Approach*, McGraw-Hill, 2000. [2](#)
- Frankel, Alexander**, “Aligned Delegation,” *American Economic Review*, 2014, *104* (1), 66–83. [9](#), [29](#)
- , “Delegating Multiple Decisions,” *American Economic Journal: Microeconomics*, 2016, *8* (4), 16–53. [9](#), [29](#)
- Guo, Yingni**, “Dynamic Delegation of Experimentation,” *American Economic Review*, 2016, *106* (8), 1969–2008. [9](#)
- Halac, Marina and Pierre Yared**, “Fiscal Rules and Discretion under Limited Enforcement,” NBER Working paper 25463, January 2019. [9](#)
- Holmstrom, Bengt**, “On Incentives and Control in Organizations.” PhD dissertation, Stanford University 1977. [1](#)
- Joskow, Paul L. and Richard Schmalensee**, “Incentive Regulation for Electric Utilities,” *Yale Journal of Regulation*, 1986, *4*, 1–49. [2](#)
- Koessler, Frederic and David Martimort**, “Optimal Delegation with Multi-dimensional Decisions,” *Journal of Economic Theory*, 2012, *147* (5), 1850–1881. [9](#), [29](#)
- Kolotilin, Anton and Andriy Zapechelnnyuk**, “Persuasion Meets Delegation,” Working Paper, February 2019. [1](#), [5.1](#), [25](#)
- Laffont, Jean-Jacques and Jean Tirole**, “Using Cost Observation to Regulate Firms,” *Journal of Political Economy*, 1986, *94* (3), 614–41. [1](#)
- **and –** , *A Theory of Incentives in Procurement and Regulation*, MIT Press: Cambridge, MA, 1993. [1](#), [2](#), [3](#)
- Loeb, Martin and Wesley A. Magat**, “A Decentralized Method for Utility Regulation,” *Journal of Law and Economics*, 1979, *22* (2), 399–404. [1](#)

- Martimort, David and Aggey Semenov**, “Continuity in mechanism design without transfers,” *Economics Letters*, 2006, *93* (2), 182–189. [9](#)
- Melumad, Nahum D. and Toshiyuki Shibano**, “Communication in Settings with No Transfers,” *The RAND Journal of Economics*, 1991, *22* (2), 173–198. [9](#), [23](#)
- Milgrom, Paul and Ilya Segal**, “Envelope Theorems for Arbitrary Choice Sets,” *Econometrica*, March 2002, *70* (2), 583–601. [30](#)
- Mylovanov, Tymofiy**, “Veto-based delegation,” *Journal of Economic Theory*, January 2008, *138* (1), 297–307. [9](#)
- Schmalensee, Richard**, “Good Regulatory Regimes,” *The RAND Journal of Economics*, 1989, *20* (3), 417–36. [2](#), [4](#)