# Mechanisms for Allocation and Decentralization of Patent Rights[*]

Hugo Hopenhayn[†]     Gerard Llobet[‡]     Matthew Mitchell[§]

January 25, 2012

## Abstract

We trace the development of the literature on the allocation of patent rights from early models of patent length through modern models of cumulative innovation by many innovators. A common useful theme is describing policies in terms of patent duration. We use one framework for considering a variety of papers in the literature, and describe how duration promises can be used as a state variable in constructing optimal allocations in dynamic problems. We tie this long standing literature into a more recent literature using a mechanism design approach to study the reward of innovation under asymmetric information. Decentralization in models of sequential innovation generates a system that incorporates self-enforcement as a requirement of the optimal policy. One interpretation of the decentralization of rights between competing innovators is as a system of mandatory buyouts paid between innovators.

# 1  Introduction

Innovation is one the centerpieces of an economy, and many generations of economists have discussed how to best promote it. Whether the allocation of market power (through a patent, for example) is necessary or not has been subject to a heated debate that dates back to, at least, Smith (1762) and Mill (1848). Schumpeter (1942) claims that monopoly rents are a necessary evil in encouraging innovation. On the contrary, Boldrin and Levine (2008) suggest that the current intellectual property rights system is not a very effective tool to generate innovation while it generates numerous distortions.

The allocation of market power through a patent is a classical way of providing incentives to innovate. It is typically argued that XV century Venice and XVII century Britain created the first two patent systems.[1] Patents have not been, however, the only way to provide incentives to innovate. Prizes awarded to firms that come up with innovations have also been a common mechanism used by governments to spur innovation.[2] Finally, governments might allocate research contracts to private firms to produce an innovation.

As discussed, for example, in Wright (1983) whether either system is optimal greatly depends on the balance between their different distortions. Market power relates the profits that firms obtain with the contribution of their innovation but at the cost of the dead-weight loss generated by an inefficient allocation. Prizes or government contracts to the extent that they make knowledge public generate no dead-weight loss but they require enormous amounts of information that are typically not available.

This review addresses two questions. First, we start by assuming that the allocation of market power rights is necessary to foster innovation and we review the literature that discusses how these rights should be optimally allocated. Second, we then turn our attention to the question of what sort of informational frictions call for rewarding innovation through monopoly profits, as opposed to the planner simply buying the inventions and placing them in the public domain (for example, through prizes). We show that environments with those information frictions generate cash payments that can be interpreted as patent fees paid to the central authority, or, in the case where many innovators work on a common project, mandatory license fees that lead to payments between innovators.

For the first half of this review we take as a starting point the classic problem of Arrow (1962) and Nordhaus (1969), allocating duration of monopoly rights in order the encourage higher quality innovations. The trade-off is the one highlighted for centuries: monopoly rights, while encouraging innovation, entail the usual sort of dead-weight loss from pricing above marginal cost.

---

[1]See Chapter 1 of Scotchmer (2006) and the references within for more details on the history of intellectual property.

[2]A typical example is the Defense Advance Research Projects Agency (DARPA), dependent of the U.S. Department of Defense. It sponsors the DARPA Grand Challenge, a prize competition oriented to create autonomous ground vehicles. Prizes, however, are not only used by governments. A recent example of its usage by private institutions is the *Google Lunar X Prize* that awards $30 million to the first privately funded team that sends a robot to the Moon. See `http://www.googlelunarxprize.org/`.

While the previous trade-off is based on the idea that there is only one innovation to implement, in practice, an important feature is that one innovation builds on the prior art; innovations are many and cumulative. In this case a new issue emerges: future innovators compete away the profits of the previous generations of innovations that enable them. One could aim to eliminate this externality by granting first innovators patent rights on all future research. However, a *hold-up* problem might arise in this case. Creators of improved products internalize that their investment will not be included in the ex-post licensing negotiations. We show that hold-up can lead to a trade-off that is interesting even when the static dead-weight loss effect is absent. Throughout, we build all of the models on a common notation and structure, so that the relationship between them is apparent.

A common theme that emerges is that a key feature of patent rights is their duration. This concept is immediate in the most classic examples of one innovator, where duration is one for one in the length of patent protection. But in more complicated models with a sequence of arrivals by heterogeneous innovators who compete for profits, patent rights evolve stochastically and are history dependent in potentially complicated ways, but in the end the concept of duration, augmented to include expectations over future events, is a key variable.

In the environments developed in the literature on the allocation of monopoly rights, the planner has complete knowledge of the innovation opportunities that lead to the products that receive the monopoly rights. In the second part of this review we relax these informational assumptions. That governments typically make decisions under incomplete information is a natural assumption. First, with complete information it is likely superior to simply pay a prize that rewards high quality research. Models with asymmetric information about the quality of the prize allow for the possibility that prizes are difficult to implement, because the prize cannot be tied to the level of quality of the innovation. Rewards through market profits, by contrast, are tightly connected to the quality of the product, and therefore provide just the incentives for innovation effort outlined in the classical literature. A system where innovations receive rewards through above-marginal cost pricing arises endogenously as a result of the information structure, lending support to the sort of model of allocation of duration.

Aside from solving the moral hazard problem of an innovator's effort at producing a high quality innovation, asymmetric information about type introduces another set of issues. First is whether or not the planner should tailor the rewards to the details of the innovator's type. In many models with complete innovation, this is immediate: if some innovations generate more output per unit of duration of monopoly rights offered, those innovations should be offered more duration. Under incomplete information, those ideas are considered in an environment where enforcement must be generated by the patent mechanism. Menus of different patent protections can be offered only if they lead to self-selection by innovators of different types. More generally, the policies studied under asymmetric innovation must sort the worthwhile innovations from the valueless: they must, again, be incentive compatible,

here in the sense that firms with innovations that should not be implemented must not choose to apply for patent protection. Thinking about enforcement issues as part of the optimal patent policy is a further advance that models of asymmetric information offer.

All of the results on sorting patents across different types of innovators are focused on the same idea: greater protection must be paid for through sometime, such as a fee, that makes the protection valuable only to the types for which the particular patent right is defined. The development follows the standard methods for sorting under asymmetric information; under a monotonicity condition, sorting can be achieved through a fee. In some cases sorting can be achieved by menus of duration versus breadth: some innovators choose longer protection that is less forceful.

In the context of sequential models the fees that implement sorting across types have a further interpretation as mandatory licensing fees. Each arrival of a new innovation dictates two payments: one to the patent authority and the other to the incumbent leader. The latter is the mandatory license fee. The former can generate different protection for different types of innovators; greater protection, in the form of a higher fee to be paid by a future innovator, requires a greater payment to the authority. Menus of patent breadth are implemented through different degrees of mandatory licensing fees. The use of mandatory licensing fees is a natural outcome, in the sense that the friction in the sequential innovation case is the hold-up that can come from patent rights. The mandatory fees address the hold-up problem directly.

The policies that borrow from the mechanism design literature on sorting all take policy to depend only on an innovator's report about the innovation it produces. Another line suggests that augmenting the reports of innovators with ex post market signals, or with signals provided by other knowledgeable parties such as competitors, could further reduce the need for costly markups as part of the reward for innovation. We finish our review by describing these approaches, and how they might be wed to the models that have been developed and studied under simple adverse selection and moral hazard.

The models described presuppose that innovation, left to its own devices, does not generate sufficient rewards. That this is the case is certainly not obvious. If knowledge spills over only with ownership, for instance through reverse engineering, then it is natural to imagine that the sale price includes those possible benefits. The analogy is to any piece of capital: the mere fact that an acorn can grow an oak tree, and in turn more acorns, does not immediately imply that an authority needs to encourage the development of acorns. However, some products may generate knowledge flows outside of purchases; simple seeing the innovation being used, or in the disclosure required to sell the innovation itself, might confer valuable innovation. Such a spillover is the most natural way to imagine that the "problem" of rewarding innovation arises.[3]

The structure of this review is as follows. In section 2 we study the complete information

---

[3]In other cases, the inefficiently might be the result of increasing returns: the technology by which ideas are produced and transmitted may not allow for positive profits at competitive pricing of ideas, as is typical in models of increasing returns.

models, from the early work of Arrow and Nordhaus through recent models of sequential innovation. In section 3 we take up the question of asymmetric information. Our discussion parallels the discussion of complete information models, beginning from static models and working up to sequential models where innovators contribute to a common project, and compete with one another.

# 2 Optimal patent duration with complete information

In much of what we do, it is useful to describe lengths of time as discounted durations. In particular, let $d$ denote the expected discounted length of time of patent protection. For instance, a patent with length $T$ units of time during which the innovator is protected has duration $d = \int_0^T e^{-rt}dt$, where $r$ is the discount rate. Note that as a result, discounting until $T$, $e^{-rT}$, can be written as $1 - rd$. The value of a payoff of one every $T$ periods is $\frac{1}{rd}$. Throughout the remainder of the discussion we set $r = 1$. This choice has no impact on any of the results, and can be thought of as a units of measure choice in the time variable.

## 2.1 Single innovation, single innovator: Optimal length and breadth

### 2.1.1 Patent length and Breadth

Early contributions to the formalization of the trade-offs in patent protection, in papers such as Arrow (1962) and Nordhaus (1969), focused on the optimal duration of this right. In these models, there is an innovator who can only profit by being given monopoly power. The size of the innovation is measured as $\Delta$, which generates profits $\pi(\Delta)$ in each instant the innovator enjoys monopoly rights; without these rights profits are competed away by imitation. If these rights are awarded for $d$ discounted periods of time, the expected discounted profits are $d\pi(\Delta)$. Given costs $c(\Delta)$ of undertaking innovation, the innovator chooses the innovation level $\Delta(d)$ that results from

$$\max_{\Delta} \ d\pi(\Delta) - c(\Delta).$$

This equation is identical to the standard incentive constraint for the innovators choice of $\Delta$. The difference from problems that we study in section 3 is, first, that it is assumed that the only instrument the planner has at its disposal is $d$, and not, for instance, prizes or fees, and second, that the cost function is assumed to be known.

Monopoly entails dead-weight losses. Let $R_M(\Delta)$ be the planner's payoff from innovation in each instant of monopoly, and $R_C(\Delta)$ be the payoff with competition. Dead-weight loss of monopoly and the existence of consumer surplus implies $R_C \geq R_M$. Denote this dead-weight loss as $L(\Delta) \equiv R_C(\Delta) - R_M(\Delta)$. The planner chooses duration to maximize

$$\max_d \ R_C(\Delta(d)) - dL(\Delta(d)) - c(\Delta(d)).$$

The trade-off here is the classic one that dates back at least to Smith (1762) and Mill (1848): more duration increases innovation, but also increases costs from dead-weight loss.

It is worth to mention a case that will be relevant as we move forward. When demand is perfectly inelastic and market power does not generate a dead-weight loss, $R_C = R_M = \pi$. In that case it is immediate that the optimal patent duration is $d = 1$; the monopolist is made the residual claimant on all returns from the innovation forever, and therefore works at the first best level of effort. This trivial outcome for the most simple model will stand in contrast to what it implies in other settings.

Another concept of a patent's degree of property rights that has received substantial attention in the literature is its breadth (or scope). Breadth may pertain to what degree of differentiation is required for firms to produce a product that competes with the patented product. The simplest idea of patent breadth is the one proposed by Gilbert and Shapiro (1990). Imagine a planner who can not only issue monopoly power, but regulate its use, with the potential to reduce both dead-weight loss and profits, by regulating the monopolist's conduct to act more competitively. Denote by $B$ the breadth of a patent. The planner gets $R_M(\Delta, B)$ when the patent is in force; it is natural to normalize zero breadth to a "valueless" patent that excludes nothing, so $R_C(\Delta) = R_M(\Delta, 0)$. As before, $L(\Delta, B) = R_M(\Delta, 0) - R_M(\Delta, B)$. Further, let profits be $\pi(\Delta, B)$, where $\pi(\Delta, 0) = 0$.

A planner chooses breadth and length to maximize

$$\max_{B,d} \; R_M(\Delta(d, B), 0) - dL(\Delta(d, B), B) - c(\Delta(d, B)),$$

where the choice size $\Delta(d, B)$ comes from the solution to innovator's problem of maximizing $d\pi(\Delta, B) - c(\Delta)$. Many papers studying the concept of patent breadth are based on particular models of competition between the patented product and copies; examples include Gallini (1992) (who focuses on costly imitation) and Klemperer (1990) (who studies a model where patent breadth determines the degree of differentiation between the patented product and competitively provided alternatives). They rely on additional structure to generate forms for $R_M$, $L$, and $\pi$.[4] Studying patent breadth in this way makes explicit the connection between breadth and exclusions: patents are exclusion rights, and the degree of that right is determined by the breadth of the patent. Models of cumulative innovation, studied below, take that connection one step further: exclusions today impact the arrival of future patentable innovations.

The model above treats the problem as one of motivating a single innovator. In a sense, it is a model where a requirement for an innovation is an "idea," which is the private property of a single potential innovator. In practice, the idea may be developed simultaneously by several potential innovators. Models of racing were pioneered by Loury (1979), Lee and Wilde (1980), Dasgupta and Stiglitz (1980), Mortensen (1982), and Reinganum (1982). While one can describe such environments using the notation developed here, we focus on the private ideas case for brevity.

---

[4]Gilbert and Shapiro (1990) analyze the case where the size of innovation is fixed. They show that in that case social welfare is convex, meaning that duration should be maximal and breadth should be as narrow as possible, subject to reimbursing fixed costs.

## 2.2 Cumulative innovation

When one innovation makes possible a future improvement or application, the notion of providing rewards through market power may generate a new conflict. When innovations and its future improvement are competitors in the same market, granting market power to one of them will be detrimental to the profitability of the other one. Thus, we have a trade-off between rewarding current innovators and stimulating future research. Resolving this conflict is central to models of cumulative innovation.

### 2.2.1 Breadth as probability of exclusion rights: Two innovators

To set the stage for the model with many cumulative innovation, first consider a model of optimal protection with two cumulative innovations based on Green and Scotchmer (1995). Consider a setup with one base innovation and one follow-up innovation (perhaps an application), which we will call an *improvement*. The improvement can be made immediately after the base innovation. Suppose that an innovation of size $\Delta_1$ is followed by an improvement of size $\Delta_2$. In this section we assume that, although their size might be uncertain, it is beyond the control of the firm, that must incur in a constant development cost $c$. The improvement is carried out by an innovator distinct to the one that achieved the first innovation.

Coasian logic dictates that there is no problem with offering broad rights to the first innovator if negotiating were free from hold-up concerns. Indeed, a broad patent together with efficient licensing allows the first innovator to maximally benefit from the innovation. Therefore it is typically optimal to award a long, broad patent to the first innovator, and allow him to extract returns from his innovation both directly and via negotiation with subsequent innovations, which infringe. This renders cumulative innovation concerns moot. The risk of hold-up is the friction used to model the potential costs of a clash of rights between innovators. To make that conflict as stark as possible, first suppose that, as in Green and Scotchmer, licensing agreements cannot be made until after both firms spend $c$. To highlight the hold-up problem, suppose that the first innovator makes a take-it-or-leave-it offers to the follower, and therefore extracts all surplus. As a result, only non-infringing improvements are made if $c > 0$, since the second innovator's outside option is zero; exclusion rights completely preclude innovation by the second innovator.

This setup highlights the motivation for studying asymmetric information in the next section. If the planner could observe all the relevant variables, the planner could focus policy on mitigating hold-up directly, by simply imposing outcomes that mirror ex ante licensing decisions. In other words, the optimal policy would simply be broad protection for the first innovator, together with mandated "licensing" outcomes to combat hold-up. We will return to this issue in the next section when we analyze the optimal mechanism in the face of asymmetric information.

Competition between innovators must lower the first innovator's profits for the rights of the innovators to come into conflict. To highlight the role of competition in lowering the first innovator's profits, suppose that these profits become zero when the follower's improvement

is allowed to be produced, so that the planner's payoff in that case is identical to the one where the initial product is produced competitively, and the improvement monopolized.[5] Further, since costs of production are assumed to be zero the payoffs to all players when the improvement is provided competitively are identical to that when both products are provided competitively. This cuts down on the set of market structures that are relevant to three: (i) one firm has rights to both products, (ii) the improvement is monopolized but its owner does not have rights to the basic innovation, and (iii) the improvement is produced competitively. We call the first case monopoly (M), the second duopoly (D) (although it includes a monopolized improvement competing with either a monopolized basic innovation, or a competitively provided base innovation), and the last competition (C). Define the planner's payoffs $R_M(\Delta_1, \Delta_2)$, $R_D(\Delta_1, \Delta_2)$, and $R_C(\Delta_1, \Delta_2)$; profits for the lead firm are $\pi_M(\Delta_1, \Delta_2)$, $\pi_D(\Delta_1, \Delta_2)$, and $\pi_C(\Delta_1, \Delta_2) = 0$. Profit for the first innovator only occur when the second firm does not innovate, and therefore are the same as $\pi_M(\Delta_1, 0)$. We similarly use $\Delta_2 = 0$ when the follow up improvement is not implemented in defining the planner's payoff. We likewise let $L_D = R_C - R_D$ and $L_M = R_C - R_M$. In order to keep things as simple as possible, we will always maintain that the planner's payoff across innovations is additive; that is, the total benefit to the planner sums the payoffs across innovations.

Since the first innovator can only profit when the improvement can be excluded, here we have two notions of duration. One is the natural counterpart of the duration discussed in the one innovation case: protection expires at some point. The other regards the size of the next innovation that will be allowed. In particular, if $\Delta_2$ is drawn from some distribution $\Phi$, the planner might choose a cutoff rule (which turns out to be optimal) so that only improvements $\Delta_2 > \hat{\Delta}$ are allowed. Thus, the duration of the monopoly rights might be written as

$$d = \Phi(\hat{\Delta})d_1,$$

where $d_1$ is the continuation duration that the planner allocates if the improvement infringes the patent. The first innovator's payoff is simply $d\pi_M(\Delta_1, 0) - c$. This pins down the needed duration $d$, since $d = c/\pi_M(\Delta_1, 0)$ is required to get the innovator to break even.

The planner needs to define three more durations: the continuation $d_1$ for the base innovator when he is not overtaken, and duration $d_2^D$ for the second innovator where he has exclusive rights to the improvement, but no exclusion rights to the prior innovation (i.e. duopoly), and $d_2^M$ units during which he has both rights to his improvement, and rights to exclude the prior innovator (i.e. complete monopoly). Since the second innovation is not sufficiently good to be implemented with probability $\Phi(\hat{\Delta}) = d/d_1$ the planner's payoff is

$$\frac{d}{d_1}(R_C(\Delta_1, 0) - d_1 L(\Delta_1, 0)) + \left(1 - \frac{d}{d_1}\right) \int_{\hat{\Delta}} (R_C(\Delta_1, \Delta_2) - d_2^D L_D(\Delta_1, \Delta_2) - d_2^M L_M(\Delta_1, \Delta_2)) d\Phi(\Delta_2)).$$

This expression can be broken into two parts. The first is a planning problem similar to the Arrow-Nordhaus case above, choosing the appropriate reward for the follower when the

---

[5]One might imagine a Bertrand game between quality differentiated products, where the higher quality improvement always sells to the entire market in equilibrium.

follower is implemented:

$$W(\Delta_1, \hat{\Delta}) = \max_{d_2^M, d_2^D} \int_{\hat{\Delta}} (R_C(\Delta_1, \Delta_2) - d_2^D L_D(\Delta_1, \Delta_2) - d_2^M L_M(\Delta_1, \Delta_2)) d\Phi(\Delta_2).$$

This choice must be made subject to reimbursing the cost of the marginal innovation, $\hat{\Delta}$, or

$$d_2^D \pi_D(\Delta_1, \hat{\Delta}) + d_2^M \pi_M(\Delta_1, \hat{\Delta}) - c = 0.$$

The problem is somewhat more complicated than the simplest model by the explicit discussion of "backwards" exclusion rights: do non-infringing patents have exclusion rights over products that came before? This question arises because we have taken explicit account of the timing of innovations. Note, however, that the linearity in the problem above implies that either $d_2^D$ or $d_2^M$ is positive; the planner simply assesses which generates less dead-weight loss per unit of profit, and uses that instrument alone, leaving a single choice variable exactly like the standard duration model.

The fundamental "clash" of rights, however, pertains to the choice of $d_1$; this problem is fundamentally different from the one discussed in the single innovator literature. We can associate, with any $d$ and $d_1$, a single $\hat{\Delta}(d, d_1)$; the trade-off between rights holders can be written as

$$\max_{d_1} \frac{d}{d_1}(R_C(\Delta_1, 0) - d_1 L(\Delta_1, 0)) + \left(1 - \frac{d}{d_1}\right) W(\Delta_1, \hat{\Delta}(d, d_1)),$$

where, recall that $d = c/\pi_M(\Delta_1, 0)$. In addition to generating the usual dead-weight loss from the prior section, the greater is $d_1$, the fewer subsequent innovations need to be excluded, since $\hat{\Delta}$ is decreasing in $d_1$. An alternative view of the clash can be seen by fixing $d_1$ and varying $d$: more rights to the first innovator, in order to make innovation more attractive, necessarily lowers the amount of subsequent innovations that can be implemented. There is a trade-off between avoiding dead-weight loss on the first innovation by keeping $d_1$ low, and implementing more second generation products through limited exclusions. This trade-off here amounts to a static breadth/length trade-off, but takes explicit account of the way in which patent breadth impacts the rate of arrival of future ideas.[6]

To see the sense in which this is different, notice that the trade-off is relevant even when the dead-weight loss from monopoly, $L$, is absent. In that case, single innovation logic would dictate giving a long-lived patent; here that conflicts with implementing future innovators. The next section considers allocations in a model with many heterogeneous innovators, but no dead-weight loss.[7]

---

[6]In a similar setup Chang (1995) shows that when courts need to determine whether the improvement infringes the original patent or not but lack information about the development cost, patent breadth should be highest when the first innovation has a very low stand-alone value.

[7]Hopenhayn et al. (2006) introduce a simple model of quality differentiation where this is the case.

## 2.3 Allocating duration with a sequence of innovations

The previous model expands the set of exclusion rights to two. First was that the innovator might be able to exclude follow-up innovations. The second is that the follower might be able to exclude prior innovations. Following O'Donoghue et al. (1998) we call these forward and backward exclusion rights.[8] In this section we focus on forward exclusion rights. Suppose that a sequence of innovations arrive, each by a different innovator. Each is an improvement on prior art; as in the prior section we assume that only the best innovation can make profits.[9] Each contributor chooses the degree of innovation, so that the contributions $\Delta$ are not fixed as they were in the prior section, but rather are a function of the reward, as in the original model of optimal patent length.

We will assume that opportunities for innovations, which we call ideas, come along periodically. With Poisson arrival rate $\lambda$, a new innovator arrives with an idea $\theta$, with atomless distribution $\Phi(\theta)$, that allows an improvement with cost $c(\Delta, \theta)$.[10] Higher $\theta$ indicates a better idea, in terms of both marginal and total cost; in other words, $c_2$ and $c_{12}$ are both negative. In this section we assume everything about timing and type of arrival is observable; in the next section we talk about interpreting and decentralizing the allocation under information frictions. It is without loss to consider the case where $\theta$ is known; it turns out that truthful revelation of $\theta$ under private information can be accomplished, under a monotonicity condition on $c$, through the use of a patent fee. The optimal allocations are unchanged.

### 2.3.1 Many innovators

First, suppose that there are many innovators, each of whom have at most one idea. When obtaining an idea, innovators are promised a period of market leadership of duration $d$. We focus exclusively on how the planner allocates scarce moments of possible monopoly power among the many innovators who arrive, since we are maintaining the assumption of no static costs of monopoly. As a result, duration is always promised to some innovator[11] If a firm is considering a level of contribution $\Delta$ it maximizes

$$d\pi(\Delta) - c(\Delta, \theta).$$

---

[8]O'Donoghue et al. study these rights in a model of vertical differentiation and a sequence of arrivals by contributors to the "state of the art" quality.

[9]It is straightforward to allow profits after the innovation is the "state of the art." Such a model is completely tractable if post-leadership profits are independent of the rights structure post-leadership; however they become intractable if made to depend on arbitrary rights structures, since the set of rights grows infinite as the set of potential rights holders grows.

[10]Riis and Shi (2012) study a class of patent policies in a sequential model with endogenous arrivals of ideas determined by free entry. They show that such an environment can lead to a concern for over-innovation and optimal statutory patent length that is finite.

[11]In the market structure introduced in Hopenhayn et al. (2006), backward rights beyond the prior innovation are irrelevant to incentives to innovate, since they simply improve profits in proportion to the other firms contribution.

We allow the planner to tailor the duration offered to the arriving idea. Here the use of duration $d$ has made the problem particularly tractable: this duration may be granted as part of a complicated history dependent optimal plan on the part of the planner. Despite that complexity, the innovator can use $d$ as a summary statistic to decide on their level of innovation, and in turn, the planner's payoff, as we assume that each arrival generates benefit $R(\Delta)$ for the planner.

The model is one where Schumpeterian forces rule, in the sense of all of the models above; a given innovator is more willing to invest the more market power he is allocated. However, as in previous models of sequential innovation, these rights clash across innovators: the innovators produce products that compete, here in an extreme form where it is assumed that only one firm can be granted any profitable rights at any point in time. The planner faces a trade-off. Promises of market power to today's innovator constrains the possible allocations that can be offered to future innovators. The planner, then, must optimally allocate these scarce instants, taking into account the heterogeneity in ideas that will occur.

We describe the planner's problem as a dynamic program. The key is to define as a state variable the commitments of duration the planner has made from prior innovators. Since only one innovator can receive monopoly profits in a given instant, the planner's implementation of future ideas is independent of how the prior commitments are distributed across prior innovators; the planner needs only to keep track of their total $d$. The planner then chooses, for any arrival of type $\theta$, the duration promise to make to the new innovator $d_2(\theta)$, and the continuation promise to prior-committed innovators, totaling $d_1(\theta)$ across all of those innovators. Define $\beta = \lambda/(1 + \lambda)$ to be the effective discount factor until the next arrival. At the time of an arrival, the planning problem is

$$V(d) = \max_{d_1(\theta), d_2(\theta)} \int (R(\Delta(\theta)) - c(\Delta(\theta), \theta) + \beta V(\tilde{d}(\theta)))\phi(\theta)d\theta, \tag{1}$$

$$s.t \quad \Delta(\theta) = \arg\max_{\hat{\Delta}} d_2(\theta)\hat{\Delta} - c(\hat{\Delta}, \theta),$$

$$\tilde{d}(\theta) = \frac{1}{\beta}(d_1(\theta) + d_2(\theta)) - \frac{1}{\lambda},$$

$$d = \int d_1(\theta)\phi(\theta)d\theta.$$

The final constraint is the key new piece of economics in the problem. It is the promise keeping constraint, that says that the planner is under obligation to provide $d$ expected discounted units of duration to prior innovators upon entering the period. The middle constraint describes the evolution of the planner's commitments: they must deliver, upon the arrival of the next innovation, the duration promises made today, reduced by the delivered monopoly power between now and the next arrival. This implies that the current duration promise must satisfy

$$d_1(\theta) + d_2(\theta) = (1 - \beta) + \beta\tilde{d}(\theta),$$

since the promise entails all the instants until the next arrival ($1 - \beta$ units of time) plus, $\tilde{d}(\theta)$ discounted by $\beta$. Rearranging for $\tilde{d}$ yields the middle constraint. The first constraint is

exactly as in the prior sections; a promise of monopoly rights encourages innovation. Given those two constraints, the planner allocates duration in the dynamic program. The planner receives benefits through $R$ depending on the award promised and the resulting level of innovation $\Delta(\theta)$.

Hopenhayn et al. (2006) show that this program implies a value function $V$ which is concave. They then develop two main results regarding the allocations resulting from the solution to (1). The first is about the duration promise that evolves, $d$. They develop a sufficient condition under which $d_2(\theta) > 0$ implies $d_1(\theta) = 0$.[12] In other words, under the sufficient condition, awarding duration to a new arrival always ends the duration promise to the previous innovator; this is termed "exclusive" rights. This makes the allocation especially simple to track: a single current "patent holder" has the duration promise $d$; a new arrival either ends the duration promise, and leads to a new patent holder, or the idea is bypassed and the prior innovator continues. In particular, exclusive policies can be described by a patent breadth $B$ such that ideas $\theta > B$ are implemented, and worse ideas are bypassed, as in the model of two innovators in section 2.2.1. In that case one can write the problem as

$$
\begin{aligned}
V(d) &= \max_{B, d_1(\theta), d_2(\theta)} \Phi(B)\beta V(\tilde{d}(\theta)) + \int_B^\infty (R(\Delta(\theta)) - c(\Delta(\theta), \theta) + \beta V(\tilde{d}(\theta)))\phi(\theta)d\theta, \\
s.t \quad & \Delta(\theta) = \arg\max_{\hat{\Delta}} d_2(\theta)\hat{\Delta} - c(\hat{\Delta}, \theta), \\
& \tilde{d}(\theta) = \frac{1}{\beta}(d_1(\theta) + d_2(\theta)) - \frac{1}{\lambda}, \\
& d = \int d_1(\theta)\phi(\theta)d\theta.
\end{aligned}
$$

It is immediate that $d_1(\theta) \equiv d_1$ is constant for $\theta < B$ since $\theta$ does not enter the first order condition for $d_1$. Further, from the first order condition for $d_1$, the choice results in $\tilde{d}(\theta)$ for $\theta < B$ such that

$$
V'(\tilde{d}(\theta)) = \eta,
$$

where $\eta$ is the Lagrange multiplier on the promise keeping constraint. From the envelope condition for $V$ it must be that $V'(d) = \eta$, which implies that $\tilde{d}(\theta) = d$ for $\theta < B$. In other words, duration is constant over the life of the patent, or, for any $d$, $d_1 = d\beta + \frac{\beta}{\lambda}$. From promise keeping, this implies that $\Phi(B(d))(d\beta + \frac{\beta}{\lambda}) = d$. This pins down, for any $d$, both $d_1$ and $B(d)$.

All that is left is to construct $d_2(\theta)$ for $\theta > B$. It is immediate that the solution to this problem is independent of $d$; if a new idea arrives and is implemented, it receives a duration independent of the past innovation's promise. The fact that implemented ideas are treated independently from the current state $d$ does not imply, however, that new arrivals, conditional on implementation, are treated equally; in general $d_2(\theta)$ may depend on $\theta$. This

---

[12]The condition relates the shape of the $c$ to the distribution $\phi$. Mitchell and Zhang (2012) study a related model where rights are shared, in the sense that both the duration of the new innovator and the old rights holder might be simultaneously positive.

leads to the second main result regarding the allocation of duration promises in Hopenhayn et al. (2006). If the marginal impact on $\Delta$ of a unit of duration is larger for better ideas, the planner takes advantage of this by assigning a greater duration promise to higher $\theta$. This implies that patents are a menu: better ideas get more protection, to take advantage of the fact that duration is more effective there. In the next section, where incomplete information about $\theta$ requires the planner to generate a patent system that is incentive compatible, the decentralization must encourage types to sort into the appropriate level of protection.

### 2.3.2 Two innovators

In the prior models of sequential innovation, each new innovation came from a distinct innovator. This made the problem of reward static along one dimension: while the planner had to think about dynamic optimization of the scarce resource of potential duration, each innovator's problem is static at the moment they develop their one innovation. Hopenhayn and Mitchell (2010) study a variant of the problem of the last section, but where all the innovations are made by two firms. They maintain the same notion of the clash of monopoly rights: only one of the two firms can be allowed to profit at any instant, since the innovators compete. However, a key difference between the case where all innovations are owned by different firms, and the model with two firms, is that an innovator may be allowed to profit from more than one of his innovations at a given point in time, since an innovator need not compete with himself.

Studying cumulative problems with a few innovators integrates the study of patent policy with the study of regulation of market structure more generally. In many industries a few firms account for a large share of innovations; it has been argued that in such industries, competition "for the market" at the innovation stage is more important than competition "in the market," the focus of conventional antitrust regulation. The sequential innovation model allows a formal structure for considering competition through innovation, and the implications for policy.

Consider an environment exactly the same as in the prior section, except that all the innovations may arrive to two innovators. An idea arrives with Poisson rate $\lambda$. Each idea is equally likely to arrive to each of the two innovators, regardless of the history of past arrivals. To keep the problem tractable, they still assume that the social planner's benefit is $R(\Delta)$, and for any promise of duration $d$ for a given innovation, the innovator chooses $\Delta$ in the same manner. For simplicity, suppose there is no heterogeneity, although one can interpret successive arrivals of many ideas to one innovator as being pieces of a large innovation.

In the previous model with many innovators, each of which contributed once, no heterogeneity and concavity of $R$ imply that every idea gets the same duration, exactly until the next idea arrives. The difference with two innovators is that the planner can make a given innovator's duration promises overlap, rewarding the innovator for more than one innovation at a given point in time. It turns out that this leads to rewards that depend on an innovators entire history of contributions.

At some abuse of notation let $R(d) = R(\Delta(d)) - c(\Delta(d))$, and assume that this is concave. First, suppose that the planner gives the current rights holder exclusive rights to *every* innovation ever invented. Since there are no static distortions, doing so is optimal since this ensures that at least one innovator is rewarded for each of their prior innovations at any point in time. Moreover, the planner allocates every instant to one innovator or another. The state is described by the promised duration $d$ to innovator one. Innovator two is implicitly promised $1 - d$. Consider the moment a new idea arrives, but without conditioning on which innovator has the idea. The planner chooses a new duration promise for innovator one $d_1$ if the idea arrives to innovator one, and $d_2$ if the idea arrives to innovator two. Unlike in the previous problem, the planner has a nontrivial choice of the allocation of the intervening time until the next arrival: $x_i \in \{0, 1\}$ takes the value one if time until the next arrival is allocated to innovator $i$, given an arrival by innovator $i$.[13] The planner's problem at the moment of a new arrival, unconditional on the identity of the firm who has the idea, is

$$V(d) = \max_{d_1,d_2,x_i} \frac{1}{2}\left[R(d_1) + \beta V(\tilde{d}(d_1, x_1))\right] + \frac{1}{2}\left[R(1 - d_2) + \beta V(\tilde{d}(d_2, x_2))\right],$$

$$s.t. \quad \tilde{d}(z, x) = z/\beta - x/\lambda,$$

$$d = \frac{1}{2}d_1 + \frac{1}{2}d_2.$$

As in the previous problem, the planner must keep its promise of $d$ units for innovator one. The last equation guarantees that he does, in expectation, by giving innovator one $d_1$ if the arrival is his (which, at the moment of arrival, is a probability $1/2$ event), and $d_2$ if the arrival is by innovator two. The equation defining $\tilde{d}$ is symmetric to the one in the previous model with many innovators; to deliver a promise of $z$ units of duration to an innovator who receives $x$ until the next arrival, the future promise must solve

$$(1 - \beta)x + \beta\tilde{d} = z$$

The value function for this case in concave and symmetric; this implies that it is maximized at $d = 1/2$. This is true because equal promises allow the planner the greatest flexibility to implement the next idea, regardless of its owner; unequal promises force the planner to treat innovators differently depending on whether they are favored or not. This is unfavorable because of concavity of $R$.

Nonetheless, Hopenhayn and Mitchell (2010) show that the duration promise indeed strays from equal promises. Intuitively, the planner trades off the cost of unequal promise against the benefit of rewarding the current arrival with a reward that encourages more innovation. For small departures from equality, the cost is not first order, but the benefit is. Formally, from the first order conditions for $d_1$ and $d_2$, we can see that $d_1 > d > d_2$; the innovator's promise rises with each of its arrivals, and therefore falls with each arrival by a competitor. This further implies that, since $R$ is increasing in $d$, duration cannot converge

---

[13]It would never be optimal to split the intervening time, so we leave out that possibility here.

to an point on the interior of $[0, 1/r]$. In other words, for a sufficient number of consecutive arrivals, an innovator is granted a promise of duration nearly equal to $1/r$. This amounts to nearly complete monopoly, at the exclusion of its competitors ideas. Optimal regulation of the industry generates incentives completely through competition for the market, and in the limit, one firm wins the market in perpetuity.

The results show that the ability of the planner to overlap rights between innovators generates a particular kind of "backloaded" rewards familiar from Acemoglu and Akcigit (2006). Once again, the key is the optimal use of scarce duration. By rewarding the innovator more after an increasing number of successes, it uses a single instant to reward an innovator for several innovations, which is a more efficient use of an instant that must be allocated to either one innovator or another. The key is that each unit of time cannot reward both innovators at once, but can reward one innovator for more than innovation that it has developed. This stands in stark contrast to the implication of this setup with a sequence of innovators, in which a sequence of identical opportunities for innovation each get the same treatment. With two innovators, the treatment of an innovation is tightly connected to the past history of the innovator who receives the idea.

# 3   Incomplete information

All the papers in the previous section bring up a common question. If monopoly power leads to the standard dead-weight losses why not using other mechanisms to foster innovation? Consider, for example, in the models discussed in the previous section, offering a cash price that equates the cost of the innovator. In that case the innovator would be indifferent between undertaking the right amount of innovation and taking any other action. This innovation could then be placed in the public domain leading to the efficient market outcome. But what if this cost or the diligence of the firm in innovation is private information? In the sequential case, if the problem is one of holdup, why not have the planner use its information to simply manage the hold-up problem directly? In this section we assume that the planner lacks complete information, and construct mechanisms in response. In the static case, these policies typically look like patents that involve fees paid on the side of the innovator. In the sequential case, the planner can implement the allocations of the previous section, but must do so through mandatory buyout fees that resemble policies that would mitigate hold-up.

The one-innovator problem resembles the classical discussion in the market regulation literature. Whereas here the planner chooses the optimal duration and subsidy that allow the firm to break even, in the previous literature the regulator needs to guarantee that by choosing the optimal combination of the use of monetary transfers and profit by above marginal cost pricing the firm can cover its total costs. Absent any other distortion the price should be set at marginal cost and the remaining cost fully reimbursed by the government. It is well-known that if cost is private information the first best allocation can also be achieved by granting the firm a reward equal to the consumer surplus that the production

generates. This is the insight of Loeb and Magat (1979). In the context of innovation, the same mechanism could be used to compute the reward that the firm should obtain. This reward should equal the consumer surplus that the innovation would generate. As before, this innovation could then be placed in the public domain.

Of course, although this is an illuminating benchmark, there are many reasons why the previous mechanism is problematic in practice. The literature on regulation has mainly emphasized the social cost of subsidizing a firm due to distortionary taxation. In that case, a planner (or a regulator) may raise the price in the market to equate the marginal distortions that the increase in the price and the taxes levied generate across markets.[14]

A second reason has been often discussed in the innovation context. Innovation is a very uncertain process, the benefits of which are difficult to ascertain, particularly for agents other than the creator of the technology. Private information on both costs and value of the innovation make the previous rewards schemes unfeasible. The only way a planner can provide incentives to undertake the right innovation is by providing the firm with some market power. In that case, a more appropriate innovation and of a larger size will tend to lead to higher profits in the market.

We study these questions in the context of the models that we have built so far: ones where the market power encourages innovation, but generates either dead-weight costs, or costs in terms of excluded future ideas.

## 3.1  Single innovation

The discussion in this section is based on Cornelli and Schankerman (1999) and Scotchmer (1999). These are the first two papers in the literature that studied how menus of patents of different duration in exchange for different fees could be used to screen different kinds of innovators. They propose to implement this menu using renewal fees.

Consider a model of a single innovation of fixed size $\Delta$, where $\pi(\Delta) = \Delta$ is the private value that the innovation generates (profits). As before, $c$ is the cost of obtaining the innovation. The patent schemes we consider correspond to monopoly for expected discounted duration $d$, in exchange of a fee $F$. Profits are

$$\Pi = d\Delta - c - F.$$

We maintain the previous notation for surplus: the innovation leads to a total flow of surplus $R_M(\Delta)$, whereas after expiration and under competition, the innovation leads to a flow of

---

[14]In the context of innovation, a consumption tax can mimic the static impact of monopoly, matching both dead-weight loss and revenue generated. It does so without any further dynamic costs that might arise due to cumulative innovation. So if the planner has access to consumption taxation, raising revenue is at least as efficient as monopoly, and possibly more so. Further, the planner may well have additional instruments of optimal taxation that lower the cost of taxation further. The problem of rewarding innovators switches from a question of industrial organization to one of public finance: how should we raise the revenue necessary to finance innovation.

consumer surplus (and also total surplus) of $R_C(\Delta)$. As usual, the dead-weight loss can be written as

$$L(\Delta) = R_C(\Delta) - R_M(\Delta).$$

Thus, the present value of social welfare is

$$W = R_C(\Delta) - dL(\Delta) - c.$$

As mentioned before, if both $c$ or $\Delta$ are exogenous and at least one of them is observable, patents will not be optimal. The optimal scheme compensates the innovator with a lump-sum payment of $F = c$, when the cost is observable, or $F = \Delta$ when only the value is observable.

Patents, however, will become optimal in two situations: when the investment of the innovator is privately observed or when both the value and the cost of the innovation are private information. We analyze these cases next.

### 3.1.1  Unobservable actions

Suppose that an invention of value $\Delta$ is the deterministic result of the expenditure in innovation by the firm. Furthermore, suppose that innovators are different in their ability to undertake research. This ability is characterized by a parameter $\theta$. Thus, we assume that an innovation of size $\Delta$ implies a cost $c(\Delta, \theta)$. We assume that $c_2 < 0$, $c_1 > 0$, $c_{11} > 0$ and $c_{12} < 0$.

An innovator facing a duration $d$ chooses the size of innovation that maximizes

$$\Pi(\theta, d, F) = \max_{\Delta} \ d\Delta - c(\Delta, \theta) - F,$$

or

$$d = c_1(\Delta(d, \theta), \theta),$$

so that $\Delta(d, \theta)$ is increasing in both its arguments. This problem is identical to the one in section 2.1.1, but specializing to $\pi(\Delta) = \Delta$. Note that, as in previous models, greater duration generates greater innovation. The new element is that innovators can choose their protection. When facing a menu of patents $\{d(\theta), F(\theta)\}$ the innovator solves

$$\max_{\hat{\theta}} \ d(\hat{\theta}) \Delta(d(\hat{\theta}), \theta) - c(\Delta(d(\hat{\theta}), \theta), \theta) - F(\hat{\theta}),$$

where the positive cross-derivative with respect to $\Delta$ and $\theta$ indicates that an incentive compatible menu of contracts is associated with longer durations to more efficient innovators.

The menu of contracts that maximizes social welfare results from

$$\max_{d(\theta), F(\theta)} \int_{\underline{\theta}}^{\infty} \left[ R_C(\Delta(d(\theta), \theta)) - d(\theta) L(\Delta(d(\theta), \theta)) - c(\Delta(d(\theta), \theta), \theta) \right] \phi(\theta) d\theta,$$

$$\text{s.t. } d(\theta') \geq d(\theta), \forall \theta' > \theta,$$

$$\Pi(\underline{\theta}, d(\underline{\theta}), F(\underline{\theta})) = 0.$$

17

As in standard mechanism design problems with monotone marginal cost, any monotone allocation (here $d(\theta)$ increasing) can be implemented through a pay-for-allocation system (here with payment $F(\theta)$); no non-monotone allocation can be supported. Therefore the planner can replace incentive compatibility with the monotonicity constraint on $d(\theta)$. A patent duration strictly increasing in $\theta$ might be optimal for two reasons. First, more duration allocated to better ideas might lead to a bigger increase in social surplus. This can be due either to a bigger impact on the size of the innovation – that is, $\Delta_{12} > 0$ – or because the social surplus from innovation is convex in $\Delta$, so that increases in the quality arising from better ideas lead to higher social surplus. Second, the cost of allocating more duration to better ideas might be small if the dead-weight loss from market power is decreasing in $\theta$.

### 3.1.2 Patent length and renewal fees

Scotchmer (1999) and Cornelli and Schankerman (1999) show that patents of different lengths can be implemented in the context of the current patent system. In the case of the latter, the optimal menu of contracts $\{d(\theta), F(\theta)\}$ can be understood as patents with different durations and different fees where a length of time $T$ is associated with a payment

$$\bar{F}(T) = F(\theta^{-1}(T))$$

where $\theta^{-1}(d)$ is the inverse of $d(\theta)$.

Although the current patent system does not allow for patents with different durations, the existence of renewal fees implies that patents with innovations of different qualities are enforced for different periods of time. If we denote the renewal fee in $T$ as $\gamma(T)$ the optimal mechanism can be implemented with renewal fees if

$$\gamma(T) = \bar{F}'(T) = \frac{F'(\theta)}{d'(\theta)} e^{-T}.$$

As a result, at time $T$ the innovator will want to renew a patent as long as the marginal gain is greater than the additional payment required. Under some regularity conditions, it is easy to see that this decision is equivalent to renewing the patent as long as the present value of future net profits is greater than the current fee.

## 3.2 Optimal breadth and length

As in the previous section, when firms hold private information, patent breadth can play an important role in the design of the optimal IP system. Hopenhayn and Mitchell (2001) considers a generic framework and provides conditions under which it is optimal to offer a menu of patents with different breadths and lengths.

In particular, consider patents defined by three dimensions, breadth $B$, statutory length $T$, and fee $F$. We assume that profits can be represented as

$$\Pi(B, T, \theta) - F - c,$$

where $c$ is the innovation cost. The parameter $\theta$ denotes an exogenous *idea* that arrives to the innovator, the value of which is the source of the private information. Furthermore, the social planner cannot observe whether the firm has obtained an idea or not, which implies that $F \geq 0$, since otherwise all firms would pretend that they have innovated.[15]

Obviously, breadth and length are assumed to increase profits for the patent holder. Social welfare, however, is decreasing in both breadth and length of the patent. In this setup, the authors show that if profits are monotonic in $\theta$ (either increasing or decreasing) and

$$\frac{\partial^2 \Pi}{\partial B \partial \theta} \geq 0, \tag{2}$$

$$\frac{\partial^2 \Pi}{\partial T \partial \theta} \geq 0, \tag{3}$$

the optimal patent menu implies zero fees. Intuitively, the planning problem is about reimbursing costs $c + F$; fees increase the necessary reimbursement. If possible, the planner always wants to substitute the use of fees to sort out types $\theta$ with sorting through a trade-off between $T$ and $B$. Doing so reduces the total amount of the costly instruments the planner must use. Under the two sorting conditions, such a trade-off is always available: some types prefer breadth, and others length.

Whether the optimal policy can take advantage of length/breadth menus greatly depends on the application. Hopenhayn and Mitchell show that the conditions are met, for example, in some cases of the horizontal differentiation model of Klemperer (1990) described in section 2.1.1. They are also satisfied in a model of innovation fertility where some patents generate competing improvements more quickly. In the latter example the intuition of the breadth/length trade-off is most clear: an innovation that generates improvements quickly leads to a greater value for breadth (because exclusions that breadth offers will happen sooner and more frequently) and less value to statutory length (since the patented product is likely to be overtaken by an improvement before the statutory length is reached).

## 3.3    Decentralization of allocations with sequential innovation

Next we address the question of optimal reward for innovation in the sequential model studied in section 2.3. We will suppose that allocations are exclusive, but that the planner can observe neither $\theta$ nor $\Delta$. This includes an inability to monitor "true" arrivals: every instant involves a pseudo-arrival of a valueless innovation. We show that, through a system of fees that can be interpreted as mandatory licensing fees, the planner can implement the allocations developed under the assumption of complete information, but where only patent rights could be used.

We allow the planner to offer a fee $F(\theta, d)$ based on a report of $\theta$ given outstanding duration $d$. Note that $F(\theta, d) \geq 0$ or the prize will be claimed by pseudo-arrivals, bankrupting

---

[15]One can interpret this assumption as stating that there is an abundance of another type for which $c$ is infinite.

the planner. In other words, prizes $F < 0$ cannot be part of the planner's policy. Conditional on truthful reporting, innovation satisfies

$$\Delta(d,\theta) = \arg\max_{\hat{\Delta}} \; d\hat{\Delta} - c(\hat{\Delta},\theta) - F(\theta,d)$$

which corresponds to the function $\Delta(d,\theta)$ used in constructing the optimal allocations under full information. As a result, if $\theta$ can be uncovered, the problem of allocating duration is identical to the one studied in the prior section. Since we show that the planner can implement the allocations that disregard the constraint on truthful reporting of $\theta$ (aside from the issue of genuine versus actual arrivals), we therefore have characterized the optimal mechanism under asymmetric information, using transfers and promises of a duration of profitability.

As in static models of asymmetric information in this section, monotonicity of the allocation guarantees that $\theta$ can be uncovered through $F(\theta,d)$. The same is true here, since the reporting problem is static. Recall that $d_2(\theta)$ is the duration offered to new arrivals; this is monotone under the assumptions of Hopenhayn et al. (2006). The optimal report $\hat{\theta}$ of an innovator maximizes

$$d_2(\hat{\theta})\Delta(d_2(\hat{\theta}),\theta) - c(\Delta(d_2(\hat{\theta}),\theta),\theta) - F(\hat{\theta},d)$$

Following typical results, any time that $d_2$ is increasing, truthful reporting can be implemented through $F$, simply by choosing the slope of $F$ such that the first order condition in the choice of reports is satisfied at $\theta = \hat{\theta}$.[16]

Hopenhayn et al. (2006) show that, when the patent system is exclusive, one can interpret the fees as a system of mandatory buyout fees. They accomplish this by decomposing $F(\theta,d)$. The key is that $F$ depends on $\theta$ and $d$ additively. This occurs because of the special structure when the optimal policy is exclusive: $d$ determines patent breadth, but for implemented innovations, their reward $d_2(\theta)$ does not depend on the prior promised duration $d$. As a result, the fee that decentralizes the allocation can be written as

$$F(\theta,d) = F_B(d) + F_S(\theta)$$

The first term $F_B$ serves to implement the appropriate patent breadth, making entry by types below $B$ unprofitable. The second term $F_S$ serves to generate truthful reporting of $\theta$ among types that are implemented, so that they can be given the appropriate reward. This part of the fee is familiar from Scotchmer (1999) and Cornelli and Schankerman (1999): innovators who are to receive greater protection must pay a higher fee, in a way that is only worthwhile for types for which the protection is intended.

One can simply interpret these fees as paid to the planner; alternatively, one can think of $F_B$ as being paid as a mandatory buyout fee to the incumbent leader.[17] One must then

---

[16]Monotonicity ensures the second order conditions are satisfied.

[17]Under this interpretation, an innovator willing to produce needs to pay a pre-specified fee to the previous market leader in order to sell in the market. This sort of compulsory licensing was first considered by Tandon (1982) in the context of one innovation.

augment $F_S(\theta)$ to account for expected future buyout fees received by the new innovator from subsequent innovations, given that the rate of those fees being paid and the amount of those fees are both dependent on the reward the new innovator receives. In particular, let that expectation of discounted buyout fees be $E(\theta)$; given the current promise is $d$ the buyout amount is $F_B(d)$ and the new innovator must also pay $F_S(\theta) + E(\theta)$ to the planner, so that the total amount paid is, net of buyout fees received, $F(\theta, d)$.

We see that the optimal mechanism is closely related to solving the underlying hold-up problem at the heart of the friction that makes patents costly. The planner, in a sense, allows hold-up to eliminate some arrivals. This is the cost of the patent system. On the other hand, the planner ensures that the patent does not eliminate sufficiently high quality arrivals by mandating a maximum licensing fee that can be charged. The planner trades off the benefits of limitations to these payments, in terms of future innovations implemented, against the cost of lower rewards for current innovators.

In the two innovator model of section 2.3.2, the lack of heterogeneity (except for whether a genuine idea has arrived or not) makes the question of decentralization simpler, but the reporting problem of the agents becomes dynamic: by misreporting an arrival of a new idea, they can alter the state in a way that may generate benefits in the future. The key to decentralization of the reporting of arrivals through fees, however, is that the payoff to an increase in duration for a given firm is always *increasing* in the firm's ability to produce innovations: increasing duration is only as valuable as the innovations that can be marketed during that duration. As a result, the agents value of duration $d$ and cumulative innovations $\Delta$ is supermodular in duration and innovation level, implying that firms with arrivals are *more* willing to pay for the duration increase that comes with an arrival than would a firm be that did not truly have an arrival. Once again, monotonicity allows truthful reporting to result from fees that increase in duration.

## 3.4   Decentralization and Patent Enforcement

The models of the previous section all consider the issue of enforcement as part of the patent design process. For instance, decentralization of the model of cumulative innovation with a sequence of innovators, the issuance of a new patent requires no arbiter: either one pays the fee, and receives the patent, or not. To be sure this is an oversimplified view of practical enforcement of patents, but all of these models, through insisting on incentive compatibility, take enforcement not as automatic, but rather as a feature the policy must address. Patent menus and patent buyouts are a set of rules for transfers to be made, and allocation of rights that are generated purely by the action of the transfer being made.

## 3.5   Uncovering private information through outside signals

The models of previous sections all used patents to the extent that prizes, although beneficial for avoiding dead-weight loss or conflict between innovators who compete, faced information

limitations. As such, the planner has an incentive to elicit additional information, if possible. Here we describe two possibilities: using market signals to determine the prize amount, and obtaining reports from another self-interested party, such as a competitor. We study both in the context of a single innovation model; we describe the use of such ideas in sequential models as a topic for future study.

### 3.5.1 Market signals

Market signals involve paying a prize based on ex-post market determined outcomes. Many of the key insights can be seen in a model with a "two point" demand structure. In particular, consider a market with $l$ agents. If the innovation is developed into a product, one agent attaches a value $h$ and $l - 1$ agents attach a value 1. Marginal cost of production is 0 and the cost of development of this fixed size innovation are distributed according to $\Phi(c)$. We assume that $h > l$, so that a patent definitely has social costs: a monopolist would price at $h$, rather than $l$, and as a result exclude consumers with valuation 1. The cost $c$, and potentially $h$ and $l$, are private information of the innovator, who chooses whether or not to innovate (i.e. pay $c$) based on that information and the reward policy.[18]

Kremer (2000) points out that prizes can, in some cases, be efficiently based on market quantities. In the environment proposed here, suppose that $h$ is known and $l$ is private information. In this case marginal cost pricing uncovers all the information about demand, and so a very simple policy emerges: price the product at marginal cost, and pay a reward to the innovator, based on sales $l$, equal to $h + l - 1$. Since the innovator gets the entire surplus in the form of the prize, it is immediate that investment decisions are first best.[19]

When Kremer's idea is extended to allow for uncertainty about *inframaginal* types, pricing above marginal cost can arise. One can interpret these policies as patents, in the sense that they involve market power. Consider the case in the two point demand setup but where $l$ is known but $h$ is unknown. This is a problem for prizes with marginal cost pricing since the first best investment requires the inframarginal valuation $h - 1$ to be transferred to the innovator, but $h$ cannot be determined from sales at marginal cost. The planner must decide, based on a report of $h$, what price to allow and what transfer to offer. The simple demand structure makes the pricing choices very simple: either price less than 1 (say, at marginal cost) or price at $h$.[20] The planner can also offer transfers depending on the report of $h$.

It is straightforward to show in such a model that the planner never offers a transfer

---

[18]The notion that the demand is only observed to the innovator can be thought of as indirectly arising from private information about a fixed innovation size $\Delta$, where different innovation sizes or types map into different demand structures.

[19]Chari et al. (2009) consider a similar situation to Kremer (2000) but with the possibility that the innovator can manipulate the market signal of $l$, for instance by purchasing his own product. If the cost of manipulation is linear in the degree of manipulation, it is immediate that only when the cost of manipulation is sufficiently large the prize system can successfully uncover $l$.

[20]As in Weyl and Tirole (2010), we do not allow mixing; as we describe below, mixing (rarely) with prices above marginal cost, together with harsh punishments if demand is unsatisfactory at the higher price, can elicit all information about demand.

for reports $h$ such that the monopoly price is allowed, since at such prices, the innovator is always extracting the full surplus the product generates at the price for which it is sold, and therefore makes optimal decisions on investment of $c$, given the resulting pricing structure. Positive transfers to patentees only serve to encourage reporting of such demand curves, for which dead-weight losses are generated. Therefore the policy can be described simply as a the transfer $\tau$ made to firms that opt for a price equal to marginal cost; that is, firms that opt for the prize. Of course, this transfer cannot be incentive compatible unless it is independent of the report, since all such types can report whichever $h$ gives the highest prize.

Given this structure, there is a cutoff type $\bar{h}$ such that firms with $h > \bar{h}$ opt for patents and firms with $h < \bar{h}$ opt for the prize, where $l + \tau = \bar{h}$. The optimal prize level trades off over-implementation against dead-weight loss. For (almost) all firms who opt for the prize, there is over-implementation: every prize taker is rewarded $l + \tau = \bar{h}$, even though only the cutoff type actually generates $\bar{h}$ surplus. However, the prize system avoids dead-weight loss.[21] Unlike Kremer (2000), prices above marginal cost arise for some innovations, in particular ones with high value for inframarginal types.

Weyl and Tirole (2010) study the much harder two-dimensional screening problem with asymmetry of information on both dimensions of demand (size of the market at price equal to marginal cost, and a parameter that represents the magnitude of inframarginal types), for a class of continuous demand functions. They allow the planner to set a price and make the reward a function of sales at that price. As such, they offer the planner a range of polices from a pure prize (pricing at marginal cost, as in Kremer (2000)), to the full monopoly price. One can interpret intermediate prices as patents with different breadth, as in Gilbert and Shapiro (1990). They show that typically the optimal policy opts for neither the complete elimination of dead-weight loss nor full implementation; envelope conditions on either end tend to push the solution toward one where the planner chooses prices that are above marginal cost, but below the monopoly price. Nonetheless, their model shares with the simple model of uncertainty about inframarginal types the feature that asymmetric information can drive optimal prices above marginal cost, i.e. a patent.

### 3.5.2 Signals from competitors

In the model of the prior section, randomization and harsh punishments can solve the problem: ask the agent for information about the entire demand curve, and nearly always choose to price at marginal cost. Rarely the planner randomizes over other prices, and imposes harsh punishments if the quantity does not line up with the report. The use of harsh punishments has several potential problems, though. First, harsh punishments may be difficult

---

[21]In a sense, such a policy is similar to Mitchell and Moro (2006), who study a planner who trades off dead-weight loss against over-transferring to a group that "loses" from elimination of the distortion generating dead-weight loss. They find a similar policy structure: for types with a high value of the inefficient transfer method (here the patent), the planner uses the inefficient transfer method. For low cost, the planner avoids dead-weight loss by compensation through a more efficient cash transfer, at the cost of overcompensating almost all types who receive the transfer.

to enforce, given limited liability constraints. Second, as Chari et al. (2009) stress, signals may be manipulable, and the incentive to manipulate may get strong when the punishments are great.

Kremer (1998) has an alternative use for randomization that does not require harsh punishments. Suppose that several other firms, in particular competitors, also know the demand curve. Then one can simply ask for their information about demand, for instance through an auction for the monopoly right; with a very high probability the monopoly right disappears, and marginal cost pricing is followed, but occasionally the auction does in fact transfer the property right to a competitor, motivating the participants to bid honestly. One can naturally see how this policy might be complementary to a "market signal" policy like the ones considered in the last section: pricing at marginal cost uncovers information about the extent of the market, an the auction uncovers information about inframarginal types, through the monopoly rents possible if the innovation's rights are owned.

Because such an approach involves other maximizing agents, rather than simply a market signal, it is natural to consider the incentives of those agents, particularly the incentives for bribery or collusion in the auction. This is especially a concern because with high probability the auction is nullified, and used only by the planner to generate some information about the patent's value.

The planner might do better, and perhaps combat collusion, with a more sophisticated mechanism that uses reports from competitors. Chari et al. (2009) use techniques from the mechanism design literature to show that, indeed, with a more complicated mechanism one can implement full information revelation at no cost, and do so in a way that makes truth telling a unique equilibrium. Even with these tools, however, there is the concern of collusion: if the innovator and the competitor can agree to make action-condition payments to one another, there is no incentive compatible way to uncover the information they share. The patentor is left to use price above marginal cost (i.e. a patent) in order to uncover information about demand, such as the inframarginal valuations in the model above.

Hopenhayn et al. (2006) consider the possibility of cross reporting between the sequence of cumulative innovators, and find that any mechanism that uses reports of one innovator about another innovation is susceptible to bribery. Considering the possibility of using reports from multiple innovators in other models of multiple innovation, for instance the model of Hopenhayn and Mitchell (2010), is a topic for future research.

# 4 Summary

Economists have been interested in the balance of the cost and benefits of rewarding innovation through market power for hundreds of years. The modern literature has formalized long held ideas and refined them to consider many aspects of the innovation process. Recent work has highlighted the fact that, even when market power is unambiguously good for encouraging a given innovation, it may have costs in terms of the generation of future

innovations. These models of cumulative innovation highlight the conflict between rights holders that produce competing products, and point to implications of this trade-off for the way in which rights are allocated.

This literature dovetails with the recent literature applying mechanism design tools to the study of patent policy. These papers push two new lines in the patent literature. First, one should think seriously about information frictions when considering the nature of optimal policy, since under complete information it is likely that a prize system dominates a costly system of monopoly rights. In order to best reward innovation, one must consider the possibilities for uncovering information both from innovators themselves, as well as from external signals of an innovation's value.

In the sequential case, under asymmetric information one can reinterpret patent fees as payments both to the patent authority, as in single innovation models, and also as buyout fees. The policy is a system of rules for transfer of rights between innovators that balances the benefits of rewarding innovation against the costs they impose for future innovators.

# References

ACEMOGLU, DARON AND UFUK AKCIGIT, "State-Dependent Intellectual Property Rights Policy," NBER Working Papers 12775, National Bureau of Economic Research, Inc, December 2006.

ARROW, K.J., "Economic Welfare and the Allocation of Resources for Invention," in National Bureau of Economic Research conference report, ed., *The Rate and Direction of Inventive Activity*, Princeton, 1962 pp. 619–25.

BOLDRIN, MICHELE AND DAVID K. LEVINE, *Against Intellectual Monopoly*, Cambridge University Press, 2008.

CHANG, HOWARD F., "Patent Scope, Antitrust Policy, and Cumulative Innovation," *RAND Journal of Economics*, Spring 1995, *26*(1), pp. 34–57.

CHARI, V.V., MIKHAIL GOLOSOV AND ALEH TSYVISKI, "Prizes and patents: using market signals to provide incentives for innovations," Working Papers 673, Federal Reserve Bank of Minneapolis, 2009.

CORNELLI, FRANCESCA AND MARK SCHANKERMAN, "Patent Renewals and R&D Incentives," *RAND Journal of Economics*, Summer 1999, *30*(2), pp. 197–213.

DASGUPTA, PARTHA AND JOSEPH STIGLITZ, "Uncertainty, Industrial Structure, and the Speed of R&D," *The Bell Journal of Economics*, 1980, *11*(1), pp. 1–28.

GALLINI, NANCY T., "Patent Policy and Costly Imitation," *RAND Journal of Economics*, Spring 1992, *23*(1), pp. 52–63.

GILBERT, RICHARD AND CARL SHAPIRO, "Optimal Patent Length and Breadth," *RAND Journal of Economics*, Spring 1990, *21*(1), pp. 106–112.

GREEN, JERRY R. AND SUZANNE SCOTCHMER, "On the Division of Profit in Sequential Innovation," *RAND Journal of Economics*, Spring 1995, *26*(1), pp. 20–33.

HOPENHAYN, HUGO, GERARD LLOBET AND MATTHEW F. MITCHELL, "Rewarding Sequential Innovators: Prizes, Patents, and Buyouts," *Journal of Political Economy*, December 2006, *114*(6), pp. 1041–1068.

HOPENHAYN, HUGO AND MATTHEW F. MITCHELL, "Optimal Patent Policy with Recurrent Innovations," 2010, mimeo.

HOPENHAYN, HUGO A AND MATTHEW F. MITCHELL, "Innovation Variety and Patent Breadth," *RAND Journal of Economics*, Spring 2001, *32*(1), pp. 152–66.

KLEMPERER, PAUL, "How Broad Should the Scope of Patent Protection Be?" *RAND Journal of Economics*, Spring 1990, *21*(1), pp. 113–130.

KREMER, MICHAEL, "Patent Buyouts: A Mechanism for Encouraging Innovation," *Quarterly Journal of Economics*, 1998, *113*(4), pp. 1137–1167.

—, "Creating a Market for New Vaccines. Part I: Rationale," in Josh Lerner Adam B, Jaffe and Scott Stern, eds., *Innovation Policy and the Economy*, volume 1, National Bureau of Economic Research, 2000 pp. 35–72.

LEE, TOM AND LOUIS L. WILDE, "Market Structure and Innovation: A Reformulation," *The Quarterly Journal of Economics*, 1980, *94*(2), pp. 429–436.

LOEB, M. AND W. MAGAT, "A Decentralized Method for Utility Regulation," *Journal of Law and Economics*, 1979, *22*, pp. 399–404.

LOURY, GLENN C., "Market Structure and Innovation," *The Quarterly Journal of Economics*, 1979, *93*(3), pp. 395–410.

MILL, JOHN STUART, *Principles of Political Economy with some of their Applications to Social Philosophy*, Longmans, Green and Co., ed. 1909, London, 1848.

MITCHELL, MATTHEW F. AND ANDREA MORO, "Persistent Distortionary Policies with Asymmetric Information," *American Economic Review*, March 2006, *96*(1), pp. 387–393.

MITCHELL, MATTHEW F. AND YUZHE ZHANG, "Shared Rights and Technological Progress," 2012.

MORTENSEN, DALE T., "Property Rights and Efficiency in Mating, Racing, and Related Games," *The American Economic Review*, 1982, *72*(5), pp. 968–979.

NORDHAUS, WILLIAM D., *Invention, Growth, and Welfare: A Theoretical Treatment of Technological Change*, Cambridge, Massachussets, 1969 .

O'DONOGHUE, TED, SUZANNE SCOTCHMER AND JACQUES-FRANOIS THISSE, "Patent Breadth, Patent Life, and the Pace of Technological Progress," *Journal of Economics & Management Strategy*, 03 1998, *7*(1), pp. 1–32.

REINGANUM, JENNIFER F., "A Dynamic Game of R and D: Patent Protection and Competitive Behavior," *Econometrica*, 1982, *50*(3), pp. 671–688.

RIIS, CHRISTIAN AND XIANWEN SHI, "Sequential Innovation and Optimal Patent Design," 2012.

SCHUMPETER, JOSEPH A., *Capitalism, Socialism, and Democracy*, Harper & Row, 1942.

SCOTCHMER, SUZANNE, "On the Optimality of the Patent Renewal System," *RAND Journal of Economics*, Summer 1999, *30*(2), pp. 181–196.

—, *Innovation and Incentives*, number 0262693437 in MIT Press Books, The MIT Press, 2006.

SMITH, ADAM, *Lectures on Jurisprudence*, University of Glasgow, 1762.

TANDON, PANKAJ, "Optimal Patents with Compulsory Licensing," *Journal of Political Economy*, June 1982, *90*(3), pp. 470–86.

WEYL, E. GLEN E. AND JEAN TIROLE, "Materialistic Genius and Market Power: Uncovering the best innovations," 2010, unpublished.

WRIGHT, BRIAN D., "The Economics of Invention Incentives: Patents, Prizes, and Research Contracts," *The American Economic Review*, 1983, *73*(4), pp. 691–707.