# The Externalities of Risky Science

Preliminary and incomplete

Michael Mandler*

Royal Holloway College, University of London

This version: April 2013

## Abstract

When scientists choose research projects with the highest expected value an externality can appear; slight variations on existing work will be selected in preference to new lines of research that would ultimately generate more value. New research enjoys two advantages: it is riskier and hence more likely to lead to high-value follow-up projects and it can generate more follow-up projects. Less communication among scientists can mitigate the externality, as can the presence of scientists who are risk-takers and secretive. When scientists maximize citations rather than expected value, the externality can but need not be alleviated.

**JEL codes:** D62, D83, H41, Z19

**Keywords:** science, innovation, externalities, citations, trees, martingales

*Address: Department of Economics, Royal Holloway College, University of London, Egham, Surrey, TW20 0EX, UK. Email: m.mandler@rhul.ac.uk.

# 1   Introduction

Thanks to the internet, scientists nowadays learn the results of their peers' research with little delay. This development might appear to be unambiguously good: armed with the knowledge of what research has already been undertaken and how successful it has turned out to be, other scientists can build on past successes and make better decisions about which projects to pursue. An externality, however, can overturn this scenario. When a researcher does not know what projects others are pursuing – or has only vague information – he or she may be forced into initiating a new and riskier line of research that has greater upside potential or that has a greater potential to generate follow-up projects. Although work on a new line of research will likely have smaller expected value than work in existing fields, if the new research turns out to be successful then other scientists can undertake a stream of related projects that will generate enough value to outweigh the costs of experimentation.

Two mechanisms are at work. First, while projects in existing fields may have higher expected value they will normally be less risky. But risk is an advantage. If a new field turns out to have high value then related projects will be available that also have high value while if the new field is of poor quality it can simply be abandoned. Since for these follower projects there is a potential upside gain but no downside loss, greater dispersion brings a net benefit. Second, even if new fields are no riskier they may generate more streams of successor projects. Consequently in the event that a new field has high value but one of its streams of successors turns out to be a failure, there will be other streams to turn to. To incorporate this second 'fallback' effect, we view new fields as trees whose branches represent follow-up projects. When either mechanism is in play, investment in a new field will be socially beneficial even though it will in expectation incur a short-term loss of research value. In some cases, therefore, impeding the flow of knowledge can lead to a social gain: scientists might experiment with new approaches when they are intellectually isolated, thus ultimately increasing the flow of research value.

The above story is motivated in part by a chapter in the history of particle physics. By the late 1960's, much of the particle physics community had rejected quantum field theory and instead followed the latest fashion, the 'bootstrap model.' Some heterogenous pockets

remained out of the loop, however, especially those distant from the epicenter of bootstrap research on the west coast of the U.S. (Ne'eman (1982)). In particular, a group of physicists in the Soviet Union – an academic 'island' in the words of Polyakov (1997) – continued to pursue a theory of gauge fields that would eventually be harnessed to describe the three fundamental forces in today's standard model of particle physics. With the triumph of the standard model, the bootstrap model faded away. The moral of the story is that it can be valuable to have several scientific schools following different lines of research in ignorance of each other's work.[1] When in contrast everyone knows exactly what every other researcher is doing and they all judge the expected value of research in the same way, then the pursuit of the highest-value projects will lead individuals to herd, with all researchers pursuing similar lines of attack.

The history of science is often portrayed as a struggle between a few heroic paradigm shifters and the larger crowd that dutifully pursues Kuhnian normal science; progress is slow in this account because of the shortage of brave geniuses. This paper will retain the Kuhn (1962) distinction between innovative and normal science but argue against the 'hero' explanation of the divide. The emphasis here will be on incentives: although it can be socially suboptimal to work in existing fields, the scientists who make this decision are not automatons blindly chasing down the routine details left behind by smarter pioneers, they pick projects that will in fact maximize the value of their own research. An echo of the 'hero' model will remain, however, in that scientists with a taste for risk and secrecy can mitigate the externality that afflicts research decisions and increase the flow of research value over time.

It might seem that if scientists maximize their citations rather than the value of their research then the externality will vanish; citation pursuers will get credit from all of the offshoots of the fields they initiate. It turns out that the pursuit of citations can but does not always encourage investment in new fields. Since a new field must show sufficiently high value for it to garner citations, the same forces that discourage value-maximizing scientists

---

[1]Ne'eman (1982) emphasizes that the development of gauge theory was spurred by the independent investigations of heterogeneous schools of physicists. Polyakov (1997) sounds a similar theme. The above account draws on the ideas of the physicist Rafael Sorkin who has stressed the advantages of limited communication among schools. For histories of this episode, consistent with the interpretation we have given, see Hoddeson et al. (1997) and 't Hooft (1999).

from innovating apply to citation-maximizers as well. For some distributions of values for new fields, citations can lead to an overshooting where excessively many new fields are initiated.

Among the many forces guiding scientific research that this paper neglects, the most prominent is the role of journal editors and referees. But we will at least see (in section 6.2) that the refereeing process can blunt the citation incentive to initiate new fields and that referees rather than editors present the larger incentive problem.

While scientists in this paper will herd into the field with the highest value projects, the logic that drives this clustering is different from the herds of Banerjee (1992) and Bikhchandani et al. (1992). It is not the case in our model that some scientists have better information and that other scientists for this reason mimic their decisions. We take the opposite tack of assuming that scientists share a common pool of information; herding is instead a consequence of the correct (symmetric information) pursuit of self-interest. A closer match is the 'learning by doing' model of Jovanovic & Nyarko (1996) where an agent can achieve long-run productivity growth only if his momentary expertise in the technology he knows best is not so great that the agent declines to experiment with new technologies. While a similar lesson holds here – access today to higher quality projects can bring about a long-term loss – the mechanics are different. In our model, riskiness is indispensable if new fields are to deliver a benefit to society whereas in J & N riskiness produces no direct social gain. New technologies in J & N instead derive their advantage from their greater long-run productivity; in our model any specific new field is a poor prospect but it is optimal to sample new fields since they can be dropped whenever they turn out to have low value.

Perhaps the work closest to the present paper is Hong & Page (2004) who argue that a population of agents who use a diverse set of problem-solving procedures can outperform a population of high-ability agents. Though the setting is different, the best-performing agents in Hong & Page suffer from the drawback that they all pursue the same solution to a problem, comparably to the scientists in our model who cluster in the same field.

Finally we mention a different inefficiency that scientific research can generate: the duplication of effort in the race to be the first to make a discovery, which amounts to another type of herding. See Dasgupta & David (1994) and Dasgupta & Maskin (1987). In the

present paper, we will assume duplication away – multiple scientists will never undertake the same project and the success of a project will only add to the expected value of its neighbors – not due to a belief that duplication is unimportant but to clarify that the externalities under discussion work by different paths. The role of secrecy in science is also discussed in Dasgupta-David but they consider its negative side not its potential to remedy inefficiency that we discuss in section 6.3.

## 2    Projects, fields, and trees of knowledge

Scientific research will proceed via a sequence of *projects* that are organized into *fields*. A project is a completed work of research ready for publication while a field is an innovation that makes a new set of projects possible. A new field's innovation can be theoretical and does not have to uncover new phenomena; if it did not carry so much freight, 'paradigm' might be a better expression than 'field'.

A new field may be initiated at any time, but within a field progress is cumulative and therefore the projects that are currently available are determined by the projects undertaken in the past. We model this relationship by assuming that each field $f$ is a *tree*: there is a root project, labeled $(0, f)$, which determines a set of successor projects or branches, each of which in turn determines a set of successor projects, and so on. Except for root projects, no project can be undertaken before its immediate predecessor has been.

For the issues we pursue, it will be enough to consider only trees where every non-root project has just one successor  So a root project initiates a field and leaves in its wake a set of successors each of which in turn has a single successor, and so on. A tree for a field can therefore be characterized by the integer number of successors, $\beta$, to the root project and accordingly is called a $\beta$-*tree*. We assume there is a $\beta > 0$ such that every field in the model is a $\beta$-tree. See Figure 1 for two sample trees. We could generalize considerably and allow uncertainty about tree structure: what is important is that root projects are expected to have more branches (successor projects) than nonroot projects. The value of projects and $\beta$ need not be related: a routine idea can have many potential applications while a deep theoretical idea may have only a few main lines of development.
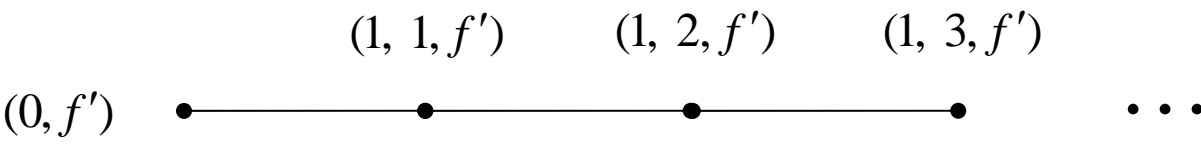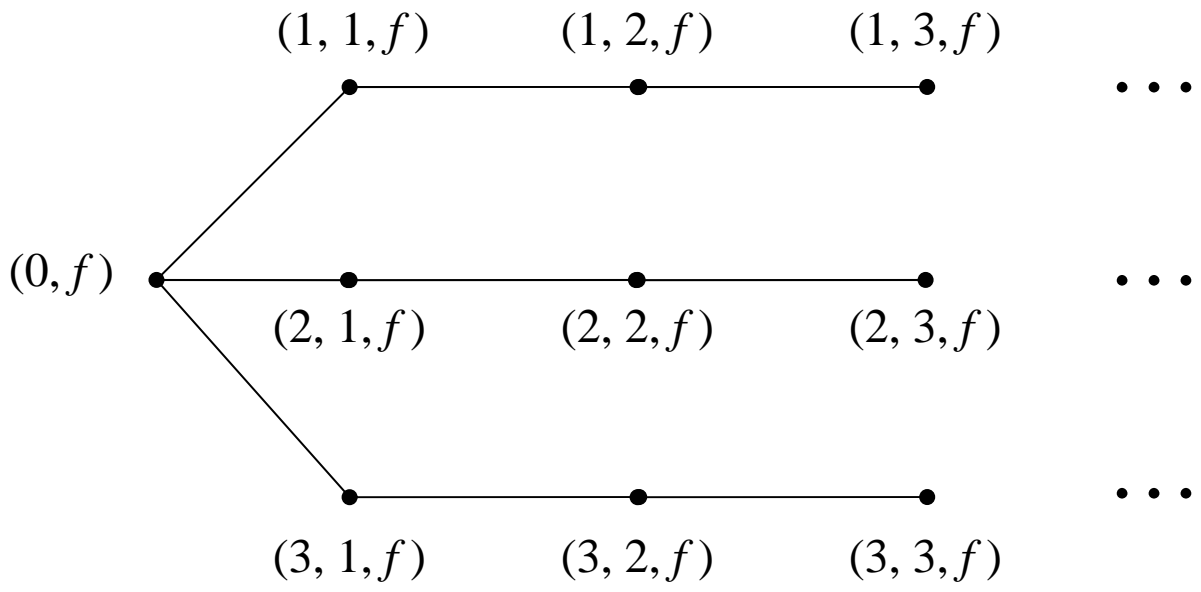
Figure 1: a 3-tree and a 1-tree

Each project in a field $f$ is identified by the triple $(b, i, f)$, where $b$ indexes the branch that leads to $(b, i, f)$ and $i$ is the number of steps between the root project and $(b, i, f)$. We use $(b, 0, f)$ as well as $(0, f)$ to denote the root project of a field $f$. So for any project $(b, i, f)$, its immediate successor is $(b, i + 1, f)$ and (if it is not a root) its immediate predecessor is $(b, i - 1, f)$. A generic project, either root or follower, is labeled $\pi$.

We consider a finite sequence of scientists $1, ..., T$ each of which undertakes a single project, where scientist $t$ chooses at date $t$. The projects that are *feasible* at $t$ are the root projects of fields that have not yet been initiated and the projects that have not yet been undertaken but that are the immediate successors of projects that have been.[2]

Each project $\pi$ has a value $v(\pi)$, a real number that indicates the project's worth and that should be interpreted as the knowledge or value added by $\pi$ given the value already generated by $\pi$'s predecessors in its field. The values of projects are uncertain and so a *state* $\omega$ will specify the value of all possible projects. $P(\cdot)$ will indicate the probabilities of sets of states.

Agents do not know $v(\pi)$ before $\pi$ is undertaken. But once $\pi$ is undertaken by scientist $t$ then $v(\pi)$ is revealed before scientist $t + 1$ decides on a project and is remembered ever after.

Let a *history* be a $h_t = (\pi_1, ..., \pi_t; v(\pi_1), ..., v(\pi_t))$ where $\pi_i$ indicates the project adopted in period $i$. For each $h_t$, we require that if $(b, i, f)$ is in $h_t$ – that is, $(b, i, f)$ is one of the first $t$ coordinates of $h_t$ – and $i \geq 1$ then $(b, i - 1, f)$ is in $h_t$, reflecting the fact that projects in a branch must be undertaken in sequence. In terms of histories, a project $(b, i, f)$ is *feasible following* $h_t$ if $(b, i, f)$ is not in $h_t$ and, when $i \geq 1$, $(b, i - 1, f)$ is in $h_t$.

A scientist's decision about which project to undertake turns on the expected values of the projects that are available. We assume that (1) ex ante any two root projects have the same expected value, which for convenience we set to equal 0, and (2) the values of projects along any branch of any field follows a random walk, i.e., the expected value of a project is equal to the value of its immediate predecessor. The upshot of these assumptions is that when a project $\pi$ in an existing field is undertaken and revealed to have positive value then

---

[2]To ensure that each scientist has the option of initiating a new field, we assume that there are at least $T$ fields.

each of $\pi$'s successors will have greater expected value than any root project $\pi'$ of a new field or any of the successors of $\pi'$. New fields will then be an unattractive source of projects. We could do without (1) but we would then have to keep track at each date of the value of the best field on offer (the uninitiated field with the highest expected value).

**Assumption 1 (common expected values for fields)** For all fields $f$ and all histories $h_t$ such that $(0, f)$ is not in $h_t$, the conditional expectation of $v(0, f)$ given $h_t$ equals 0: $E(v(0, f)|h_t) = 0$.[3]

The reader should not assume that 0 or negative value projects are worthless or destructive. Even a line of research that fails on its own terms imparts knowledge, if only the knowledge that a candidate theory is false or, in technological applications, that an invention does not work.

For (2), we assume that the expected value of any nonroot project equals the value of the project's predecessor. Given a history $h_t$, let a *leading project* of $h_t$ be a project $(b, i, f)$ in $h_t$ such that $(b, i + 1, f)$ is not in $h_t$.

**Assumption 2 (random walk)** For any history $h_t$ and any leading project $(b, i, f)$ of $h_t$,

$$E(v(b, i + 1, f)|h_t) = v(b, i, f).$$

Even though random walks or martingales represent fair gambles, we will see that the possibility of dropping negative-value fields will make it possible to achieve a positive flow of value.

The most prominent model that satisfies Assumptions 1 and 2 is a *simple random walk*, where any root project either has value 1 or $-1$, each with probability $\frac{1}{2}$, while any other project has the value of its immediate predecessor plus or minus 1, each with probability $\frac{1}{2}$. Whenever possible we will use simple random walks to make our main points, though typically the steps of the random walks will vary from field to field, rather than always equaling 1.

---

[3]Here and subsequently $E(x)$ will be the expectation of the random variable $x$ and $E(x|A)$ will be the conditional expectation of $x$ given the event $A$. We define $h_t$ as an event explicitly in Appendix A.

It is easy to introduce diminishing expected returns to the pursuit of a line of research by replacing Assumption 2 with $E(v(b, i+1, f)|h_t) = \delta v(b, i, f)$, where $\delta < 1$. Then, if a project along some branch has positive value the expected value of its successors would revert to the mean, which under our normalization is 0.

Since projects with negative expected value play little role in our analysis, we could let them deviate from Assumption 2 and assume instead that a successor to a project $(b, i, f)$ with $v(b, i, f) < 0$ has an arbitrary negative expected value. This generalization would allow a branch or field to fail with no chance of recovery, e.g., when its capacity to yield further insight is completely exhausted. Our structural assumption that any project has a potential successor would then have no bite: the endless supply of projects in a field could all be nearly worthless (have a highly negative value).

Assumption 2 implicitly rules out any correlation of expected values across branches of a single field: once the value of some root project $(0, f)$ has been discovered then every immediate successor of $(0, f)$ that has not yet been undertaken has the conditional expected value $v(0, f)$, regardless of what $v(b, 1, f)$'s have been revealed. But the correlations that are ruled out could go either way. For example, suppose $(0, f)$ is a theoretical innovation and that each immediate successor amounts to an application of this original insight. It could well be that when one application, say $(b, 1, f)$, is unusually successful – that is, $v(b, 1, f) > v(0, f)$ – then one should infer that the theory works in the real world and hence that a second application $(b', 1, f)$ will be successful too. Equally it could be that all of the applications are minor variations on a theme; hence only the first application delivers a substantial incremental insight, with the remainder delivering nearly the same message. In the first case $v(b, 1, f)$ and $v(b', 1, f)$ will be positively correlated and in the second they will be negatively correlated. Assumption 2 steers a middle course.

Given Assumptions 1 and 2, the law of iterated expectation implies that conditional on a history $h_t$ *any* project in a field that has not yet been initiated in $h_t$ has 0 expected value. In contrast any successor of a leading project $\pi$ with $v(\pi) > 0$ has positive expected value. Despite this fact, we will see that the expected value produced through time can be increased by dropping fields with positive-value leading projects and instead initiating new fields.

## 2.1 Technological change

Our model intentionally mimics Kuznets's celebrated (1930) account of technological innovation and thus indirectly draws on Schumpeter (1911) too. In Kuznets, every sector in an economy eventually falls victim to the law of diminishing returns to technical innovation. But though each sector is doomed to eventual stagnation, the economy as a whole can grow robustly due to the neverending supply of new 'leading' sectors. A corresponding pattern appears in the present model. Due to the random walk assumption, any successful field whose current projects have high value will eventually run through its stock of great ideas: given enough time, the value of projects along any branch will turn negative. But although any one new field ex ante generates projects of 0 expected value, the opportunity to switch to new fields can lead to an endless stream of projects with positive expected value.

While we have laid out the model in reference to scientific projects, the borrowing from Kuznets indicates that the model can be applied to technological innovation. Consider a large set of firms that all produce a good subject to rapid technological change. Each firm takes the good's price and prices in the future to be given exogenously, say because the firms reside in an open economy that produces for the world market, and maximizes expected profits. In each period, one firm builds a new factory and an investment decision determines the cost of output at the factory. Firms observe the technologies used by their predecessors and can copy any of them exactly, giving them access to any cost level that has been achieved in the past. But a firm can also innovate and adopt a new technology with an uncertain cost. Some innovations are entirely new; we suppose there is a large supply of such untried innovations and that each will lead to a decrease in expected costs equal to some modest nonnegative level (perhaps 0). But if a firm is lucky some of its predecessors could have discovered a major process invention that can be further refined, leading to an expected fall in costs that might be substantial.

It should be clear that this example is a special case of our model of science. Each firm is a scientist, an investment in a new factory is the undertaking of a project, the reduction in costs relative to factories that have been built in the past corresponds to the value of the project, the uncertainty of the cost reduction is the uncertainty of project value, an untried

innovation is the initiation of a new field, and the refinement of an existing innovation is a project along a branch of an existing field. It need not be that a 0 value for projects in our model means no reduction in costs: a 0 value for a project can be interpreted as the default level of per-period cost reduction possible in this industry.

There may appear to be one discrepancy between the two settings. We have assumed that scientists cannot simply repeat projects that have been previously undertaken whereas a firm can adopt an existing technology. But we could have allowed scientists to repeat previously undertaken projects. Since the value of a project is its *addition* to knowledge, the value of repetition would be the minimal negative value permitted by the model (putting considerations of plagiarism aside). A scientist seeking to add to knowledge would therefore never choose to replicate a project. A firm is in the same position: given the availability of some expected cost reductions, it will not choose to exactly replicate an existing factory.

Given our assumption that firms can freely copy any existing factory design, it is not surprising that externalities are present. Firms might not undertake an untried innovation since they will not reap all of the gains if the innovation turns out to be successful. But in this purely economic setting, a cure for the externality suggests itself: give firms the right to sell or license the right to further develop the technology it has invested in. We will exploit this suggestion in section 6.1, where we lay out a market solution for the externalities that appear in the creation of scientific knowledge.

One advantage of the application to technological change is the assumption of price-taking, which disposes of the monopoly effect that a technological innovation has on the price of its output. Since there is no reduction of innovation due to a monopoly effect, the model pinpoints the pure externality induced by the opportunity to copy other firms' inventions. In addition, all of the welfare effects of technological change are due to the accumulation of profits (i.e., cost reductions): there is no division of benefits between firms and consumers.

# 3    Plan and equilibrium

A *plan* consists of $T$ functions $a = (a_1, ..., a_T)$ where each $a_t$ assigns to each history $h_{t-1}$ a project that is feasible following $h_{t-1}$. Given the probabilities of states, a plan defines a probability that any given set of projects is undertaken and an expectation of the sum of the value of the scientific research undertaken from 1 to $T$. From society's point of view, the success of science is measured by the magnitude of this expected value. The exact definitions of these probabilities and expectations are tedious and so we segregate them into Appendix A.

We assume initially that each scientist maximizes the expected value of the project that he/she undertakes.

**Definition 1** *Given the history $h_{t-1}$, a project $\pi$ is an* equilibrium choice *(at $t$) if*

- *$\pi$ is feasible following $h_{t-1}$,*

- *$E(v(\pi)|h_{t-1}) \geq E(v(\pi')|h_{t-1})$ for all $\pi'$ that are feasible following $h_{t-1}$.*

*An* equilibrium *is a plan $(a_1, ..., a_T)$ such that each $a_t$ assigns an equilibrium choice to each history $h_{t-1}$.*

Assumptions 1-2 imply that to identify an equilibrium choice when facing $h_{t-1}$ the scientist at $t$ need only look at the leading projects of $h_{t-1}$, choose the successor of a leading project with the highest value if that value is positive, and otherwise initiate an arbitrary new field.

# 4    A simple random walk

To see how an equilibrium proceeds and illustrate the Kuznets/Schumpeter implications of the model, suppose the value of projects in each field follows the simple random walk introduced in section 2. Formally, for each $h_t$,

$$P\Big(v(0,f) = 1 \,\Big|\, h_t\Big) = P\Big(v(0,f) = -1 \,\Big|\, h_t\Big) = \frac{1}{2} \tag{4.1}$$

if $(0, f)$ is not in $h_t$ and

$$P\Big(v(b, i, f) = v(b, i-1, f) + 1 \,\Big|\, h_t\Big) = P\Big(v(b, i, f) = v(b, i-1, f) - 1 \,\Big|\, h_t\Big) = \frac{1}{2} \quad (4.2)$$

if $(b, i-1, f)$ is a leading project in $h_t$.

Since fields are all ex ante identical, let the date 1 scientist initiate an arbitrary field $f$. Suppose fields have just one branch ($\beta = 1$) and – in contrast to what happens in equilibrium – that all scientists at dates $t > 1$ simply undertake the sole available project in $f$ that is currently available. Since there is then no further initiation of new fields after date 1 the sequence of realized values will form a simple random walk and hence the expectation (at date 0) of the value of research achieved at each date $t$ will be 0.

An equilibrium performs better than this benchmark. Still assuming that each field has a single branch, the scientists at dates $t > 1$ will undertake the immediate successor of the project $\pi_{t-1}$ undertaken at $t - 1$ if $v(\pi_{t-1}) > 0$ and will initiate a new field if $v(\pi_{t-1}) < 0$. (If $v(\pi_{t-1}) = 0$ the date $t$ scientist can either select $\pi_{t-1}$'s successor or initiate; assume for concreteness that $t$ chooses $\pi_{t-1}$'s successor.) This strategy is the mirror image of the classic gambling strategy of continuing to place bets until one's stake hits positive territory: scientists in equilibrium pursue a field until its value turns negative. Now in any equilibrium the expected value of the *terminal project* of a field $f$ – the project in $f$ that is chosen immediately prior to the initiation of another field or at period $T$ – will equal 0, unsurprisingly since project values in a field in effect form a series of fair gambles.[4] This conclusion will hold in any model meeting the random-walk Assumption 2, not just in simple random walks. In Figure 2, the values of terminal projects when $T = 4$ are recorded at the terminal nodes; each fork in the figure represents the two possible values that a project might have, not a set of successors to a project (we have set $\beta = 1$). Since in equilibrium a new field is initiated when an existing field hits the value $-1$, a new field whose expected value in the subsequent period is 0 must be inserted at those points, as pictured in Figure 3. This gain from replacing an expected value of $-1$ with 0 implies that the expected value of research in each period $t > 1$ must be strictly positive. In the first four periods of any equilibrium,

---

[4]Since the values of projects in a single field form a martingale, this conclusion follows from the Doob stopping theorem (see Williams (1991)).

```
                                    0
                        /                       \
                    -1                           1
                                            /         \
                                        0               2
                                      /   \           /   \
                                    -1     1         1     3
                                          / \       / \   / \
                                         0   2     0   2 2   4
```
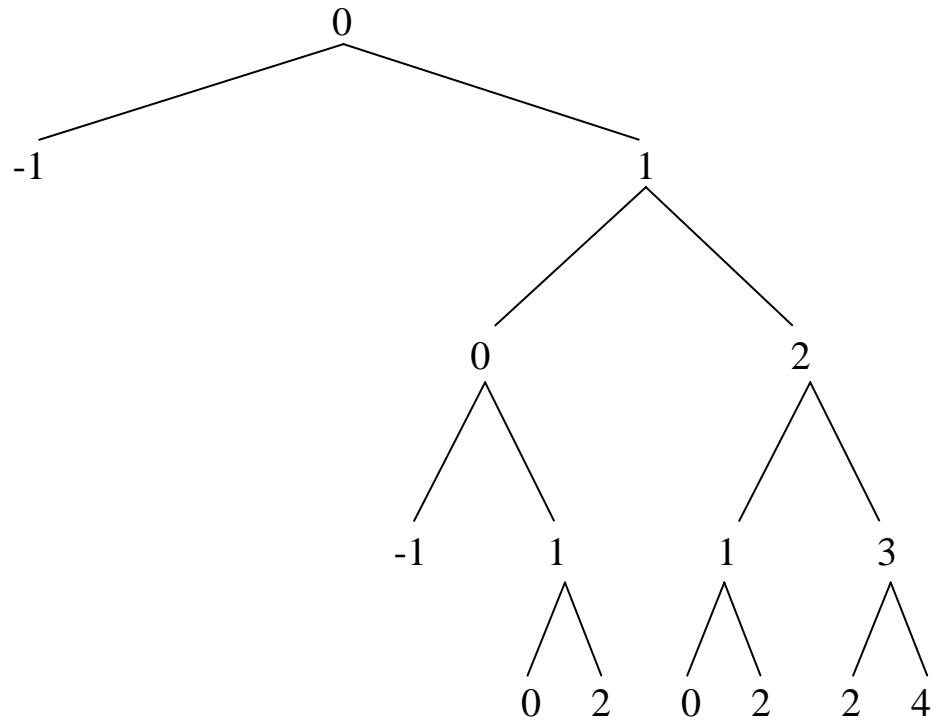
Figure 2: possible project values in a single field

```
                                    0                          Expected Value
                        /                       \
                    -1                           1                    0
                  /       \                   /       \
                -1         1               0             2           1/2
              /   \       /   \          /   \         /   \
            -1     1     0     2       -1     1       1     3        3/4
           / \   / \   / \   / \      / \   / \     / \   / \
         -1   1 0   2 -1   1 1   3  -1   1 0   2   0   2 2   4        1
```
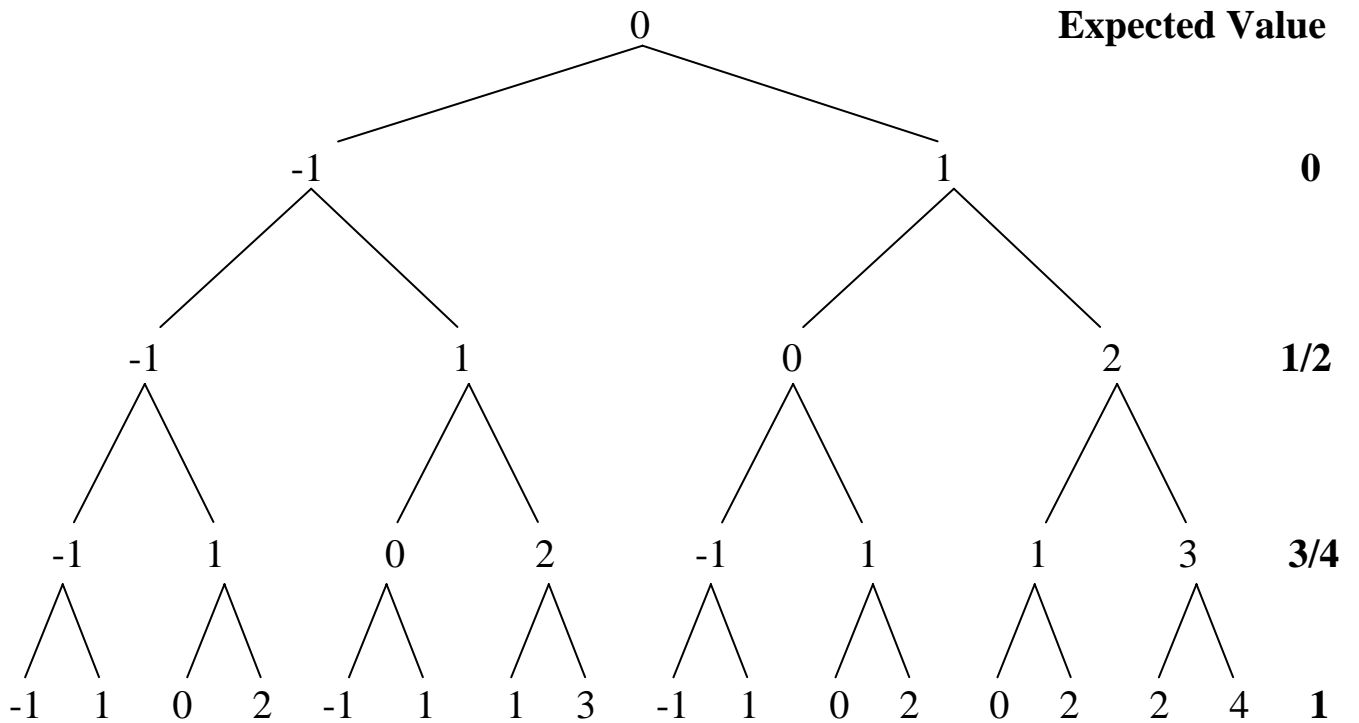
Figure 3: possible project values with new field initiation

the expected values of research turn out to be $0, \frac{1}{2}, \frac{3}{4}, 1$.

The curiosity of our random-walk assumption is that the expected value of the terminal project in any newly initiated field remains exactly equal to 0. It is only due to the new fields that are initiated whenever the projects in an existing field turn negative that the stream of expected values becomes positive overall.

It so happens that in any equilibrium for this example the expected value of research achieved at date $t$ increases without bound as $t$ increases, a consequence of the fact that the value of projects in a field does not diminish in expectation as more are undertaken. If we had followed the diminishing returns version of Assumption 2, discussed in section 2, where for some $\delta < 1$

$$P\Big(v(b,i,f) = \delta v(b,i-1,f) + 1 \,\Big|\, h_t\Big) = P\Big(v(b,i,f) = \delta v(b,i-1,f) - 1 \,\Big|\, h_t\Big) = \frac{1}{2}$$

rather than (4.2), then the expected value of research per period would be bounded above. But equilibria would still achieve a stream of positive expected values: even in the presence of diminishing returns in every field, the expected value of research in equilibrium will be strictly positive at each $t > 1$, comparably to the Kuznets growth model.

## 5   The externality

Society's interest lies in maximizing the total expected value of research, defined formally in Appendix A. An equilibrium of self-interested scientists, each maximizing the value of his or her own project, can fail to achieve this goal; in particular, the initiation of new fields can fall short of the optimal level.

If some leading project $\pi$ in an existing field has positive value then every project in any new field has smaller expected value than any successor of $\pi$ (see section 2). If science were risk-free, that would be the end of the story: society will be better off if scientists undertake the successors of the positive-value projects in existing fields and that is what self-seeking scientists will do. But we will see in the next subsection that the greater riskiness of new fields can mean there is a social gain to investing in them, even when positive-value fields

are available. We then apply these conclusions to show that there can be advantages to scientists working in isolation. In subsection 5.3, we turn to the second benefit of new fields: when a new field has multiple branches and one branch turns out to be a failure, there are fall-back branches to turn to.

## 5.1 The social benefit of risk

In the simple random walk of section 4, a social gain is achieved in equilibrium when a field whose available projects have a negative expected value is abandoned and replaced by new fields. The magnitude of this gain is given by the difference between the abandoned expected value $-1$ and the new expected value $0$. Evidently the magnitude of this gain would be greater if the dispersion of values around a project's expected value were greater than 1: then the gain achieved in the event that the values in the field turned negative would be even larger. This effect could be large enough that it might be efficient to abandon an existing low-dispersion field in favor of a new field even though the old field's projects still have positive expected value. Of course, a self-interested scientist would always choose one of the remaining positive expected-value projects in the low-dispersion field; the root project of any new field has an expected value of 0 and the scientist would not take into account the expected social gain that can occur in a new field's later periods of operation.

We therefore augment the simple random walk of section 4 by letting the riskiness of scientific research vary by field. Suppose each field $f$ is characterized by a *value increment* $s_f > 0$: the possible values of project $(0, f)$ are $-s_f$ and $s_f$ rather than 1 and $-1$, each with probability $\frac{1}{2}$, and each subsequent project $(b, i, f)$ can have value $v(b, i-1, f) + s_f$ or $v(b, i-1, f) - s_f$, also with probability $\frac{1}{2}$ each.

Value increments will be uncertain, however, until a field $f$ has been initiated when the value $v(0, f)$ will reveal $s_f$. We assume that each field is ex ante identical and more specifically that each $s_f$ is drawn independently from a common distribution with a continuous density $g(s)$ defined on the nonnegative numbers with $g(0) > 0$.[5] We call the model of this section, a simple random walk with uncertain value increments and an arbitrary fixed

---

[5]Formally, if $(0, f)$ is not in $h_t$ and $A \subset \mathbb{R}_+$ is measurable then $P(v(0, f) \in A | h_t) = \frac{1}{2} \int_A g(s) ds$ and if $(b, i-1, f)$ is a leading project of $h_t$ then in (4.2) 1 and $-1$ are replaced by $|v(0, f)|$ and $-|v(0, f)|$.

number of branches $\beta$, a *risky random walk*. For future reference, in a *risky random walk with diminishing returns*, the two possible value realizations for $(b, i, f)$ are $\delta v(b, i-1, f) + s_f$ and $\delta v(b, i-1, f) - s_f$, where $\delta < 1$.

We will say that a plan $a$ *underinvests in new fields following history $h_{t-1}$* if $a$ does not initiate a new field following $h_{t-1}$ and there is an alternative plan $a'$ that is identical to $a$ through $t-1$, that does initiate a new field following $h_{t-1}$, and that increases the total expected value of research conditional on $h_{t-1}$. See Appendix A for more details.

**Proposition 1** *In a risky random walk there is a cutoff $c > 0$ for the value of projects such that, for all $t < T$ and all $h_{t-1}$ where the value of all leading projects lies in $(0, c)$, any equilibrium underinvests in new fields following $h_{t-1}$. There is consequently a positive probability that an equilibrium will underinvest in new fields at some history.*

**Proof.** Let $w > 0$ be the maximum value of the leading projects at $h_{t-1}$. A new field $f$ initiated following $h_{t-1}$ then leads to a sacrifice in expected value in period $t$ of $w$. If $s$ is the value increment of $f$ then there is an expected benefit in period $t+1$ equal to $\frac{1}{2}(s - w)$ if $s > w$ and $0$ if $s \leq w$. Hence the expected benefit in periods $t$ and $t+1$ of initiating a new field is greater than the cost if

$$\frac{1}{2}\left(\int_w^\infty (s - w)g(s)ds\right) > w.$$

Since $\frac{1}{2}\left(\int_w^\infty (s - w)g(s)ds\right) - w$ is a continuous function of $w$ and equals $\frac{1}{2}\int_0^\infty sg(s)ds > 0$ when $w = 0$ we may set $c$ so that the above inequality is satisfied for all $w \in (0, c)$. As $c \to 0$, the maximum total value from $t+2$ to $T$ that can be generated from a field whose leading project at $h_{t-1}$ has a value $w \in (0, c)$ converges to $0$ and hence for small enough $c$ cannot overtake the gain at $t$ and $t+1$ from initiating a new field. ∎

Proposition 1 holds for any fixed $\beta$, including $\beta = 1$. In the $\beta = 1$ case, the only differences between a new, untried field $f_n$ and an existing field $f_e$ whose leading project has a positive value is that (1) the root project of $f_n$ has lower expected value that the project currently available in $f_e$ and (2) the value increment of $f_n$ may differ from that of $f_e$. Effect (1) argues in favor of the project in $f_e$ that will in fact be taken in equilibrium.

14

For equilibrium behavior to be suboptimal, therefore, it must be that there is an advantage to the possible risk characteristics of new fields. Since in a risky random walk the fields with low-value leading projects must also be low risk, Proposition 1 implies that it is the high-risk fields that deliver the greater benefit. The calculation in the proof of Proposition 1 shows that if existing fields have low value and hence low risk then a new field is likely to have higher risk and hence lead to a larger value gain when a low realization materializes.

The social advantage of risky fields is both a surprising and general conclusion. Given any plan that undertakes a project $\pi$, if we replace $\pi$ with a project $\pi_r$ such that the distribution of $v(\pi_r)$ is a mean-preserving spread of the distribution of $v(\pi)$ (and shift the means of the distributions of the successors of $\pi_r$ to preserve the random walk assumption) then there is a plan that uses $\pi_r$ instead of $\pi$ that produces a weakly larger expected value of research. The reason is simply that, following some history $h_t$, a plan that undertakes $\pi$ and a set of $\pi$'s successors for a given sequence of outcomes for $\pi$ and the projects undertaken after $\pi$ can be replaced by plan where $\pi_r$ and its successors are undertaken in the same order when the same sequence of outcomes obtains. Since, at $h_t$, $\pi_r$ and each of its successors has the same expected value as $\pi$, the new plan must generate the same expected value as the old plan. However, due to the mean-preserving spread, one of the successors of $\pi_r$ may turn out to have such a small value that it would increase total expected value to instead undertake a project in a different branch or field. Using $\pi_r$ rather than $\pi$ can therefore strictly increase expected value.

## 5.2 No news can be good news: the advantages of research is isolation

We can use Proposition 1 to show that the expected value of scientific research can sometimes increase when scientists work in isolation, cut off from detailed knowledge of the work of their peers. We stress the word 'sometimes'. When scientists work in isolation then clearly it could be beneficial, both for the scientists and for society at large, for the scientists to learn the results of their peers' research: some scientists might be pursuing projects of little value while their peers are achieving great things. But the reverse case, that communication can

sometimes be welfare diminishing, is also possible.

Suppose to begin that the world scientists is divided into $n \geq 2$ separate 'schools', each of which works in isolation. These schools are an idealization of the different camps of particle physicists described in the Introduction. The scientists in any one school are aware of which projects are undertaken by other schools but either do not know the value they have delivered or do not know enough about the projects to undertake their successors. Since as before we assume that each scientist aims to maximize the additional value generated by his or her scientific research, scientists will not repeat the projects that have been already undertaken by other schools.

The schools will choose projects in a fixed rotation from periods 1 through $T > 0$. In the first $t > 0$ of these $T$ periods the schools operate in isolation. So if $j$ is one of these periods then, given the history $h_{j-1}$, we redefine a project $(b, i, f)$ to be *feasible at $j$* if $(b, i, f)$ is feasible as defined in section 2 and if in addition, when $i \geq 1$, $(b, i-1, f)$ is chosen in history $h_{j-1}$ by the same school that chooses at $j$. Following period $t$, full communication reigns and we revert to our prior definition of feasibility. An equilibrium plan remains as stated in Definition 1 except that the amended conditions on feasibility are assumed to hold. If there is no shift to a full-communication regime, then we say that 'isolation continues'.

When full communication begins at $t$, the immediate response is for new field initiation to fall relative to what would happen if isolation were to continue. This initial effect holds with probability 1 for any model where it is negligible event for a feasible nonroot project to have an expected value exactly equal to 0 and lasts for at least $n$ periods. To see what drives this decline, notice that if isolation were to continue then, with probability 1, a new field will be initiated in one of the next $n$ periods, say $i$, only if all of the projects feasible under isolation for the school that chooses at $i$ have nonpositive expected value. If, however, full-communication begins at $t$ then the school that chooses at $i$ will have the additional option of choosing projects from fields initiated by other schools and these projects may have strictly positive expected value. So the chooser at $i$ under full-communication will initiate a new field only if every feasible project in *every* field has nonpositive expected value. Consequently, if there are $K$ schools that would choose from existing fields from $t$ through $t + n - 1$ if isolation were to continue then there must be at least $K$ projects that

must be undertaken once full communication begins before *any* school will initiate.[6]

**Proposition 2** *Suppose it is a 0-probability event for a feasible nonroot project to have an expected value equal to 0. Then in the first n periods after full communication begins, the number of new fields initiated will with probability 1 be no greater than the number of new fields that would be initiated in the same n periods if isolation were to continue. With positive probability the number of new fields initiated will be strictly smaller in the first n periods after full communication begins.*

**Proof.** The remaining 'positive probability' claim follows, for example, from the fact that it is a positive probability event for each school prior to the beginning of full communication at $t$ to select only root projects that turn out to have negative value under isolation and for the school that chooses at $t$ (under isolation and under full communication) to select a positive-value root project, say in field $f$. Under full communication the school that chooses at $t+1$ will then choose from $f$, while if isolation were to continue every school from $t+1$ to $t+n-1$ would initiate. ∎

Proposition 2 does *not* imply that full communication will reduce the expected value of scientific research: if some schools have discovered high expected-value lines of research then the initiation of a new field that an isolated school might undertake could well be counterproductive. But it is certainly possible that the positive expected-value projects that other schools have to offer are not worth pursuing. Indeed the moral of Proposition 1 is that if projects in existing fields have low but positive expected-value and therefore display low risk then the initiation of new fields will increase the subsequent expected value of research. If this condition on expected values in existing fields obtains, then the expected value of research will increase if isolation continues for another period and the school that chooses next is forced to innovate.

---

[6]The reasoning we have given applies with probability 1 and not to every state due to the knife-edge possibility that a country could decide under autarchy to undertake a project in a field that the country has already initiated even though that project has 0 expected value but, when full communication begins, instead decide to initiate a new field (assuming every feasible project has 0 expected value).

**Proposition 3** *In a risky random walk, for any date $t \leq T - 1$ sufficiently large and any equilibrium under isolation there is a positive probability set $H$ of histories from 1 to $t - 1$ such that, given any $h_{t-1} \in H$,*

- *if isolation were to continue during period $t$ then a new field will be initiated at $t$,*

- *if full communication begins at $t$ then a new field will not be initiated at $t$,*

- *the conditional expected value of research from $t$ to $T$ given $h_{t-1}$ will be greater if isolation continues after $t$ than if full communication begins at $t$.*

The Proposition is easily proved: it is a positive probability event that through date $t - 1$ every school except the school $c$ that chooses at $t$ draws a root project with a positive value near 0 and that $c$ (if it has already had a chance to choose a project) has drawn only negative-value root projects. The expected value of research from $t$ to $T$ will be higher if $c$ initiates a new field at $t$, as it will under isolation but will not under full communication.

## 5.3 The underexploitation of multi-branch fields

Suppose in the simple random walk of section 4 that the number of branches $\beta$ is greater than 1 and that the scientist who chooses at date 1 selects a field $f$ whose root project turns out to be successful, $v(0, f) = 1$. Subsequent scientists will select projects from some branch, say branch 1, and it will be an optimizing choice for each individual to continue choosing from this branch until $v(1, i, f)$ hits 0, as it will eventually with probability 1. The scientists who follow will then undertake projects in another branch of $f$. The values in this branch have expected value 1 rather than the expected value of 0 that would obtain in a new field. Multi-branch fields thus offer a plain advantage over single-branch fields.

But equilibria do not fully exploit the opportunities offered by multi-branch fields. If a scientist has a choice between a project with expected value 0 in an existing field and initiation of a new field, he will be indifferent between the options. The initiation of a new field is, however, superior from society's point of view. If a new field $f'$ turns out to have the high realization $v(0, f') = 1$ and the first successor of $f'$ that is undertaken turns out to have value 0 – which is the low realization given that $v(0, f') = 1$ – then the next

scientist can adopt another of $(0, f')$'s immediate successors, which will have expected value 1. When the same chain of realizations occurs with a 0-expected-value project in an existing field – first a high realization then a low realization – there is no 'fall-back' option available with expected value greater than 0. So far, the externality is not very impressive since it disappears if scientists choose to initiate a new field rather than select an existing-field project with expected value 0. But the externality becomes robust when we leave the narrow confines of a simple random walk.

To pinpoint how a multiplicity of branches can by itself lead to a suboptimal failure to initiate new fields, we need to neutralize variations in the risk characteristics of fields – otherwise the source of benefit in initiating a new field might just be the potential of a new field to display greater risk, which we have already analyzed. We therefore now assume that any two feasible projects have the same distribution of values after correcting for the difference in their means. Formally, there will be a density $g$ with an interval support such that if a project $\pi$ is feasible following history $h_t$ and $E(v(\pi)|h_t) = e$ then the conditional distribution of $v(\pi)$ given $h_t$ is governed by the density $g_e$ defined by $g_e(v) = g(v-e)$.[7] When this assumption holds (in addition to Assumptions 1 and 2), we say that *fixed dispersion* is satisfied.

Under fixed dispersion, the distribution of values of a root project of a new field has a less advantageous upper tail than the upper tail of any project in an existing field with positive expected value: for any $r > 0$, any $h_t$, any feasible $\pi$ such that $E(v(\pi)|h_t) > 0$, and any $f'$ that has not yet been initiated,

$$P\Big(v(\pi) > r \,\Big|\, h_t\Big) > P\Big(v(0, f') > r \,\Big|\, h_t\Big).$$

Yet under fixed dispersion it remains possible that an equilibrium underinvests in new fields. It is here that the tree structure of fields comes into play. When fields have more than one branch, the initiation of a new field has a signal advantage: if the root project $(0, f)$ has high value but the first successor of $(0, f)$ undertaken turns out to have an unexpectedly low value then other scientists can pursue the higher expected-value projects available in

---

[7] So $\int_{-\infty}^{\infty} v g_e(v) dv = e$ and, for any measurable $A$, $P(v(\pi) \in A|h_t) = \int_A g_e(v) dv$.

the remaining $\beta - 1$ branches of $f$. In a single branch of an existing field or in a 1-tree, in contrast, there are no such fall-back options. Consequently, when $\beta \geq 2$ and the expected values of projects in existing fields are sufficiently low it will be worthwhile from society's point of view to sacrifice those small expected values to gain the fall-back advantages offered by a new field.

**Cutoff Lemma** *If fixed dispersion holds then the following two conditions are equivalent:*
*(1) there exists a cutoff $c > 0$ for the value of projects such that, for all $t \leq T - 2$ and all $h_{t-1}$ such that the highest value leading project has a value in $(0, c)$, any equilibrium plan underinvests in new fields following $h_{t-1}$,*
*(2) the number of branches in each field is greater than or equal to 2.*

Since for any $c > 0$ there is a positive probability that the project undertaken at date 1 turns out to have a value in $(0, c)$, the Cutoff Lemma implies that underinvestment at date 2 will occur with positive probability if fields have at least two branches and there are at least two additional periods following date 2. The Cutoff Lemma also leads to the converse that if fields have just a single branch then equilibria never underinvest in new fields under fixed-dispersion.

**Proposition 4** *If fixed dispersion holds then with positive probability an equilibrium will underinvest in new fields at some date if and only if $T \geq 4$ and the number of branches in fields is greater than or equal to 2.*

# 6    Remedies for the externality, complete and partial

The remainder of the paper considers various solutions to the externalities discussed in section 5. The first, suggested by the application to investments in technology in section 2.1, is to set up markets for the right to work on projects. Markets overcome the externalities under consideration by rewarding agents if they undertake risky projects or multibranch fields that prove successful: the right to work on the successors of their research will then sell for a high price. Of course, markets for the right to follow up on scientific discoveries do not exist and would carry many drawbacks. They would for example run counter to

the ethos of an unfettered pursuit of knowledge, which helps drive the intellectual culture of science. But there are advantages to markets that suggest some features of the other remedies we consider. A market for the right to do research in a field is forward-looking: the returns to innovation that would normally accrue to the agents who act only in the future can instead be channeled to the innovators. Scientists who pursue recognition of their work as measured by the citations they receive are also forward-looking in the sense that if their research is successful they will earn a flow of future rewards (and can anticipate that flow as soon as their work achieves present success). After citation maximization, we look at another technique for capturing the external benefits of a scientific success, keeping research secret.

## 6.1   A market for projects

We assume that each scientist $i$ maximizes the expectation of the sum of his consumption and the value of the project undertaken in period $i$. Agent $i$ is initially endowed with ownership rights for some (possibly empty) set of root projects. The owner of any project, root or nonroot, can sell it to other agents or undertake the project and sell its successors. Any agent can also package together any set of feasible projects that he owns to sell as a bundle. The purchase of a bundle imparts the right to undertake any feasible project $\pi$ in the bundle, to sell the successors of $\pi$, and to resell the projects that are not undertaken. Packages can be unbundled at will.

Bundling is necessary for markets to achieve efficiency. To see why, consider a simple model with just two periods and where $\beta \geq 2$. There are two types of fields: fields where the value of every project equals 0 and fields that are simple random walks (they have value increments of 1). In the first period, the efficient choice is to pick a root project $\pi$ of one of the random walks. But if the successor projects of $\pi$ were not bundled, then the market price of each must equal 0 even if $\pi$ is successful and has value 1: the successors of $\pi$ will be in excess supply since only one period remains and $\beta \geq 2$. Bundling readily solves the problem since the chooser at period 1 can bundle the successors of $\pi$ into one package and sell the package for price 1 in the event that $\pi$ is successful.

To formalize, we assume that each feasible project is owned at each point in time by

some agent. Let $H$ be the set of histories. Following an arbitrary $h \in H$, agent $i$ chooses a partition $\mathcal{P}_h^i$ of the feasible projects that he owns into a set of *bundles*. If together the $T$ agents form $n_h$ bundles, then the endowment of $i$ following $h$ can be represented by a $n_h$-vector $e_h^i$, each coordinate equaling either 1 for ownership or 0 for non-ownership. A purchase by $i$ of a subset of bundles is given by a $n_h$-vector $b_h^i$, again consisting of 1's and 0's.

Given the bundling decisions of the agents following history $h$, and letting $p_h$ be the $n_h$-vector of prices for bundles, agent $i$ faces the budget constraint

$$x_h^i + p_h \cdot b_h^i \leq p_h \cdot e_h^i, \tag{BC}$$

where $x_h^i$ is agent $i$'s consumption following $h$. If $i$ undertakes a project ($h$ contains $i-1$ periods) then the undertaken project must be in one of the bundles that $i$ buys ($b_h^i$ must equal 1 in that bundle's coordinate). Let $h'$ be one of the immediate successors of $h$. The endowment of $i$ at $h'$ then consists of the projects that $i$ owns (which have a 1 entry in $b_h^i$) but does not undertake and the successors of $\pi_h$ in history $h'$ if $i$ undertook $\pi_h$ following $h$. Each agent then makes a new bundling decision following $h'$, and so forth.

A *price sequence* $p$ specifies a price vector for every history $h$ and every partition of the feasible projects at $h$. Letting $H^{i-1}$ denote the histories that contain $i-1$ periods, a *plan* for agent $i$ is a $\left( (\mathcal{P}_h^i, b_h^i, x_h^i)_{h \in H}, (\pi_h)_{h \in H^{i-1}} \right)$. Given the plans for the other agents, a plan for $i$ is *budget feasible* if BC is satisfied for all $h \in H$, $i$ buys $\pi_h$ for all $h \in H^{i-1}$, and the rules for the partitioning of endowments are satisfied.

A $\left( (\mathcal{P}_h^i, b_h^i, x_h^i)_{h \in H}, (\pi_h)_{h \in H^{i-1}} \right)$ for each $i = 1, ..., T$ determines a probability distribution over $H^t$ for each date $t$. histories at each date. Given the distribution over $H^{t-1}$, let $x^i(t)$ and $\pi(t)$ respectively denote the random variables equal to $x_h^i$ and $\pi_h$ at $h \in H^{t-1}$.

**Definition 2** *An equilibrium is a $p$ and a $\left( (\mathcal{P}_h^i, b_h^i, x_h^i)_{h \in H}, (\pi_h)_{h \in H^{i-1}} \right)$ for each $i = 1, ..., T$ such that*

- *each $\left( (\mathcal{P}_h^i, b_h^i, x_h^i)_{h \in H}, (\pi_h)_{h \in H^{i-1}} \right)$ is budget feasible and the $Ev(\pi(i)) + \sum_{t=1}^{T} Ex_h^i(t)$ achieved by $\left( (\mathcal{P}_h^i, b_h^i, x_h^i)_{h \in H}, (\pi_h)_{h \in H^{i-1}} \right)$ is at least as great as any budget-feasible alternative,*

- *for each $h$,*

$$\sum_{i=1}^{T} b_h^i \le \sum_{i=1}^{T} e_h^i,$$

  *and where $p_h(k) = 0$ if the above inequality is strict in coordinate $k$.*

Unfortunately there can be equilibria in which agents make inefficient decisions. For the simplest example, suppose $\beta = 1$ and there are at least three periods  The agent in the second-to-last period can either choose (a) a project from a field with certain value $w > 0$ whose successors are also sure to have value $w$ or (b) a project from a field that is a simple random walk (with value increment 1).  The efficient decision is to choose project (b) if $w < \frac{1}{2}$.  But suppose in the event that project (b) is successful that equilibrium prices in final period for projects (a) and the successor of (b) are 0 and $\frac{w}{2}$ respectively (exploiting the fact that both the supply and demand for the successor of (b) will equal 1 at any price between 0 and $w$).  If (a) is chosen, let final-period prices be $w$ and 0.

Given these anticipated prices and date $T - 1$ prices of $w$ for (a) and 0 for (b), agent $T-1$ will select project (a) which earns a profit of $w$ rather than (b) which earns an expected profit of $\frac{w}{4}$ (project (b) earns 0 expected value at $T - 1$ and its successor sells for $\frac{w}{2}$ with probability $\frac{1}{2}$).

But there is always an efficient equilibrium.

**Proposition 5** *There exists an equilibrium such that the projects chosen maximize the total expected value of research.*

## 6.2   The citations game

As we have seen, when scientists care only about the value of their own projects they ignore the riskiness and multiple-branch advantages of new fields; the result is that too few fields are initiated.  Scientists who pursue citations in some cases can help to overcome the multi-branch externality: since the initiator of a new field $f$ will get a citation credit from the projects undertaken in every branch of $f$, scientists will see at least some gain in initiating a field rather concentrating narrowly on the immediate value delivered by the project they undertake.  On the other hand, since a new field must prove to be sufficiently valuable

for it to earn citations, citation seekers can sometimes make the same decisions as the value-maximizing scientists we have considered so far. Scientists can also be swayed too much by the lure of the citations that can be earned from a multiple-branch field. A citation-maximizer may well innovate when it would be more productive instead to pick the highest-value project.

Suppose that scientists cite all of the work that made their choice of project $\pi$ possible – all of the projects in the same field that are predecessors of $\pi$ – and that each scientist seeks to maximize the number of projects that cite his or her work. More precisely, agents $1, ..., T-1$ will maximize their expected number of citations and agent $T$ will maximize the expected value of research undertaken in period $T$. Nothing in our analysis would change if instead all agents were to maximize a weighted sum of their expected number of citations and the expected value of the research they personally undertake as long as the weight on citations is sufficiently large. We will now call the agents and equilibria of section 3 *value-seeking*. The pursuit of citations introduces a strategic dimension that did not appear in the value-seeking model; scientist $t$'s choice of project will now be shaped by $t$'s expectations of which projects future scientists will cite.

The values of projects will evolve according to the risky random walk model described in section 5.1. To make it easier for citation maximization to overcome the externality that leads value-seekers to fail adopt a socially beneficial innovation, suppose that $\beta$, the number of branches per field, is large. Since each branch offers another set of potential citations to a field originator, the incentive for innovation is then heightened. As we saw in section 5.1, value-seeking scientists can suboptimally fail to innovate when the expected values of feasible projects in existing fields are positive but very small. Consider the case where only one branch in one field $f$ has an available positive expected-value project $\pi$, where $\pi$ has expected value $s_f$ (the value increment of $f$) and $s_f$ is itself small. To see that citation-maximizers can innovate in this scenario, suppose to the contrary that all agents are citation-maximizers but always choose the highest expected value project available. To calculate the expected number of citations earned by initiating a new field $f'$, notice that since the maximum expected value of projects in existing fields is the small number $s_f$, the probability that the root project of $f'$ will have a value that outstrips the expected

24

value of projects in existing fields will equal nearly $\frac{1}{2}$. At least the next $\beta$ scientists will therefore choose projects from $f'$. So the expected number of citations is near to $\frac{1}{2}\beta$ or higher. The expected number of citations earned from pursuing project $\pi$ in the existing field $f$ will in contrast be near 0 in the .5 probability event that $\pi$ turns out to have value $0$.[8] Since there are many other chains of value realizations that will lead the successors of $\pi$ eventually to be abandoned, the decision to initiate the new field $f'$ will earn a higher expected number of citations than $\pi$ when $\beta$ is sufficiently high. It thus cannot be the case that citation-maximizers always choose the highest expected value projects available.

Given the plan $a = (a_1, ..., a_T)$, agent $t$ earns the *citation* $\pi'$ at state $\omega$ if (1) $\pi'$ is a direct or indirect successor of the project $\pi$ that agent $t$ selects at $\omega$ and (2) some agent that chooses after $t$ selects $\pi'$ at $\omega$.[9] A plan $a$ thus defines a total number of citations $C_t[a](\omega)$ for $t$ at $\omega$ and hence an expected number of citations $E(C_t[a])$. In the definition below, we use $a^{\langle \pi, h_{t-1} \rangle}$ to denote the plan that coincides with $a$ in every coordinate except that $a_t^{\langle \pi, h_{t-1} \rangle}(h_{t-1}) = \{\pi\}$.

**Definition 3** *A citations equilibrium is a plan $a = (a_1, ..., a_T)$ such that, for each agent $t$ and history $h_{t-1}$, $a_t(h_{t-1})$ is feasible following $h_{t-1}$, and*

- *for each $t = 1, ..., T-1$, $E(C_t[a]|h_{t-1}) \geq E(C_t[a^{\langle \pi', h_{t-1} \rangle}]|h_{t-1})$ for all $\pi'$ that are feasible following $h_{t-1}$,*

- *for each $h_{T-1}$, $E(v(a_T(h_{T-1}))|h_{T-1}) \geq E(v(\pi')|h_{T-1})$ for all $\pi'$ that are feasible following $h_{T-1}$.*

By viewing the $a_t$ functions as strategies, a citations equilibrium qualifies as a Nash equilibrium; that an equilibrium assigns an optimizing choice to each history amounts to a subgame perfection requirement.

We begin by clarifying that citations address only the multi-branch externality. In a world of 1-trees, citation maximizers can take the same actions that value-maximizers take.

---

[8]The expected number of citations is not exactly 0 since $\pi$'s successors will be among the large set of projects with expected value 0 and some of these could in principle be undertaken.

[9]These requirements can be stated as (1) $\omega \in [\pi]_t$ and (2) $\omega \in [\pi']_i$ for some $i > t$.

**Proposition 6** *In a risky-random walk with 1-trees, any value-seeking equilibrium is a citations equilibrium.*

When the number of branches $\beta$ is sufficiently large, the actions of citation maximizers can deviate significantly from value-seeking behavior. In the proof of the Proposition below, the key argument is that when the expected values of feasible projects is sufficiently low then the value of the root project of a new field will with substantial probability overtake existing project values and therefore earn numerous citations.

**Proposition 7** *In a risky random walk with $\beta \geq 4$ and $T \geq 5$, for any citations equilibrium there exists a positive-probability set of histories at which an agent will undertake a project with an expected value less than the maximum available.*

Citation maximizers can sometimes initiate new fields when value-maximizers fail to, but it is less clear if or when this additional initiation is socially beneficial. The incentive for a citation-maximizer to initiate a new field is greatest when $\beta$ is large and projects in existing fields have low value since then a new field has a good chance to outstrip the available expected values and, if it does, enjoy a large set of followers. Since Proposition 1 showed that in a risky random walk it is in fact optimal to initiate a new field when projects in existing fields have sufficiently low value, citation maximizers can in fact increase efficiency, i.e., raise total expected value.

To pin down behavior when an agent faces a choice among multiple projects with an expected value of 0, we will now assume that citation maximizers in this situation initiate a new field and call a citations equilibrium with this property *plausible*. Initiation is the natural choice for a citation maximizer in the face of expected values tied at 0, since new fields have more descendants. But we need an explicit assumption to avoid self-feeding equilibria where agents always return to the same branch of a field whenever it offers a project with 0 expected value or even undertake negative expected-value projects based on the expectation that later agents will undertake its descendants. Another reasonable approach is to assume that agents randomize when facing ties, with equal weight on all projects that offer the same expected value. Propositions 8 and 9 would then still hold.

**Proposition 8** *In a risky random walk with $\beta \geq 4$ and $T \geq 5$, for any plausible citations equilibrium there is a positive-probability set of histories $H$ at which a new field is initiated even though some projects in existing fields have positive expected value. For any $h \in H$, if instead agents make value-seeking choices at $h$ and all continuations of $h$ then the conditional total expected value of research given $h$ would decrease.*

Unfortunately citation maximization can also diminish total expected value. Specifically, in a simple random walk or a risky random walk that approximates a simple random walk, it is never optimal to initiate a new field when positive expected-value projects are available. Rather than paying the penalty of undertaking a project with an expected value less than the maximum available, it is better to wait until existing project values have all hit zero value: there is no cost in deferring new field initiation until necessary and there is always a chance that new field initiation will not be necessary after all.

To define what it means to 'approximate' a simple random walk, we define a sequence of risky random walks given by $\langle g_n \rangle$ to *converge to a simple random walk* if the distributions defined by the $g_n$ converge in distribution to the distribution that assigns probability 1 to the value increment 1. This sequence is fixed below: by 'all risky random walks sufficiently near to a simple random walk' we simply mean 'for all $n$ sufficiently large'.

**Proposition 9** *For all risky random walks sufficiently near to a simple random walk, if $\beta \geq 3$ and $T \geq \beta + 3$ then for any plausible citations equilibrium there is a positive probability set of histories $H$ at which new fields are initiated such that, for any $h \in H$, if agents instead make value-seeking choices at $h$ and all continuations of $h$ then the conditional total expected value of research given $h$ would increase.*

If we step outside of the confines of the model and let fields have a variety of different $\beta$'s, $\delta$'s (the factor at which expected values decrease when diminishing returns are present), and expected values, the inefficiency of citation maximization becomes transparent: citation maximizers will avoid fields whose root projects have enormous expected value but tiny $\delta$'s in favor of fields with lower initial expected value but high $\delta$'s.

We have supposed that the results of any project will spread to the entire scientific community. But if journals oversee the dissemination of research and thus the flow of citations

then scientists will take on only those projects that journals will publish. Unfortunately the root project of a new field will not cite any past work and hence no citation-seeking editor or referee will have an incentive to let the research pass through the gate. The path by which citations potentially can mitigate underinvestment in new fields may therefore be blocked. There are countervailing forces: editors may derive a citation-like credit from stewarding a journal that initiates new fields. For anonymous citation-seeking referees, however, there is no cost to snuffing out a new field.

## 6.3   Secrecy partially offsets the externality

We have assumed so far that each scientist chooses at a single date. The assumption has not yet carried much significance. If we let the diminishing return version of Assumption 2 hold, then when a scientist chooses multiple projects at dates that are far apart the selection at earlier dates will have only a slight impact on the value of available choices later on. But if a scientist chooses at multiple dates and can keep the results of his or her research secret – say by delaying publication – then the scientist can initiate a new field and, when field's root project has great value, reap the rewards of its high-value follow-up projects. The externalities of risky science are then alleviated.

A scientist who undertakes a single project will never initiate a new field as long as some existing field has projects with positive value. Suppose now that a scientist can conduct $\tau$ projects in secrecy and that other scientists will not know enough of the details of these projects to undertake any of their successor projects until the scientist finally releases his research.[10]   For concreteness, let the model be a risky random walk with diminishing returns (see section 5.1). The scientist undertaking the secret research could then well prefer to initiate a new field. This conclusion follows from the proof of Proposition 1: if existing fields have small enough but positive values then an ability to keep even one project secret will be sufficient to induce a value-maximizer to initiate a new field.

Similar results apply to the fixed-dispersion model. We can reverse the original purpose

---

[10]Other scientists might of course repeat the entire sequence of projects that the scientist under considera-tion has undertaken in secret. But when the sequence begins with the initation of a new field our assumption of an ample supply of new fields means that this scenario is remote. Also, if other scientists know which projects are being undertaken in secret they will not want to replicate them.

of the Cutoff Lemma to conclude that if $\beta \geq 2$ and $\tau \geq 3$ then a scientist who can keep secrets will be better off initiating a new field if the value of projects in existing fields is sufficiently small.

Of course a scientist working in secret will not internalize the whole of the externality considered in section 5. As long as $T$, the time span of the entire model, is greater than $\tau$, the number of projects a scientist can secretly undertake, scientists will continue to suboptimally ignore some of the socially beneficial consequences of new field initiation. But the present analysis does suggest a more generous view of scientists who cagily refuse to discuss their work; even if motivated by paranoia their secrecy could well foster the initiation of new lines of inquiry, which is a socially productive goal.

# 7   Conclusion

The individual pursuit of scientific value – or the pursuit of the rewards that accompany successful scientific careers – does not necessarily maximize the total value produced by the entire community of scientists. Even when scientists seek the recognition of other researchers in the form of citations, they may avoid innovative projects that could make a rich supply of follow-up projects available. On the other hand, a taste for risk and secrecy – not traits that academia normally encourages – can ease the externality.

Our analysis has been geared to scientific research but it applies to any pursuit where individual projects or works can be ranked as better or worse, and where projects in a specific area build on the work done earlier. We mentioned the technological interpretation of the model in section 2. The trees we have used to link a scientific project to earlier work could also describe the bridge between past and present in cultural and artistic endeavors. And individuals in these fields can also pursue a citations-like credit for the work they stimulate.

# Appendix A: technical definitions

To define the probability that a specific project is undertaken in a given period, let $v(\pi, \omega)$ denote the value realized for project $\pi$ in state $\omega$, let $h_0$ denote the null history that agents face in period 1, and let $\Omega$ denote the entire set of states. Given plan $a = (a_1, ..., a_T)$, the

set of states where project $\pi$ is undertaken at date $t$, which we write as $[\pi]_t$, and the set of states where $\pi$ is undertaken at date $t$ and realizes value $v$, written $[\pi, v]_t$, are defined by

$$
\begin{aligned}
[\pi]_1 &= \{\omega : a_1(h_0) = \pi\} \text{ (equal to either } \varnothing \text{ or } \Omega) \\
[\pi, v]_1 &= \{\omega : a_1(h_0) = \pi \text{ and } v(\pi, \omega) = v\} \\
[\pi]_2 &= \{\omega : \exists v_1, \pi_1 \text{ such that } a_1(h_0) = \pi_1, v(\pi_1, \omega) = v_1, a_2(\pi_1; v_1) = \pi\} \\
[\pi, v]_2 &= \{\omega : \exists v_1, \pi_1 \text{ such that } a_1(h_0) = \pi_1, v(\pi_1, \omega) = v_1, a_2(\pi_1; v_1)) = \pi, v(\pi_2, \omega) = v\} \\
&\ \ \vdots \\
[\pi]_t &= \{\omega : \exists (v_i, \pi_i)_{i=1}^{t-1} \text{ such that } a_i((\pi_j)_{j=1}^{i-1}; (v(\pi_j, \omega))_{j=1}^{i-1}) = \pi_i \text{ and } v(\pi_i, \omega) = v_i \\
&\qquad\qquad\qquad \text{for } i = 1, ..., t-1 \text{ and } a_i((\pi_j)_{j=1}^{t-1}; (v(\pi_j, \omega))_{j=1}^{t-1}) = \pi\} \\
[\pi, v]_t &= \{\omega : \exists (v_i, \pi_i)_{i=1}^{t-1} \text{ such that } a_i((\pi_j)_{j=1}^{i-1}; (v(\pi_j, \omega))_{j=1}^{i-1}) = \pi_i, v(\pi_i, \omega) = v_i \text{ for} \\
&\qquad\qquad i = 1, ..., t-1, \text{ and } a_i((\pi_j)_{j=1}^{t-1}; (v(\pi_j, \omega))_{j=1}^{t-1}) = \pi, v(\pi, \omega) = v\}.
\end{aligned}
$$

Although our notation will not indicate the dependence, keep in mind that the events $[\pi]_t$ and $[\pi, v]_t$ are always defined relative to a plan $a$.

Probabilities are defined as usual from the relevant sets of states; for example the probability that project $\pi$ is undertaken at date $t$ when plan $(a_1, ..., a_T)$ is adopted is $P([\pi]_t)$. We define *total expected value of research* achieved by $a$ to be

$$
\int_\Omega \sum_{[\pi, v]_i : \omega \in [\pi, v]_i} v \, dP(\omega).
$$

To remove any trace of ambiguity: given a $\omega \in \Omega$ the summation above is taken over all $[\pi, v]_i$ such that $\omega \in [\pi, v]_i$ for some $i \in \{1, ..., T\}$, project $\pi$, and $v \in \mathbb{R}$. Since $\omega$ specifies a value for each project and there are only finitely many projects that can be undertaken by period $T$, there are only finitely many such $[\pi, v]_i$ and hence the summation is well-defined.

We can also define the event where the history $h_t = (\pi_1, ..., \pi_t; v(\pi_1), ..., v(\pi_t))$ occurs given $a$ by $[h_t] = \bigcap_{i=1,...,t}[\pi_i, v(\pi_i)]_i$, and accordingly the probability of a set of histories. Given $a$, we say that a state $\omega \in \Omega$ *leads to the history* $h_t$ if $\omega \in [h_t]$.

We define an *equilibrium $a$ to underinvest in new fields following $h_{t-1}$* if there exists an alternative plan $a'$ such that, for all $\omega$ that lead to $h_{t-1}$, $a_t$ selects a successor of a leading project following $h_{t-1}$ while $a'$ (1) is identical to $a$ up to $t-1$, (2) initiates a new field following $h_{t-1}$, and (3) increases the conditional total expected value of research given $\omega$.

When the particular $t$ at which underinvestment occurs is immaterial, we say that *an equilibrium $a$ underinvests in new fields at some date given $\omega$* if there is an alternative plan $a'$ such that (i) for all $t$ and all $h_{t-1}$, if $a$ initiates a new field following $h_{t-1}$ then so does $a'$ and (ii) $a'$ achieves a greater conditional total expected value of research than $a$ given $\omega$.[11] Details aside, the primary way $a'$ can do better in expectation than $a$ is to initiate new fields at histories where $a$ does not.

To consider the consequences of letting the time horizon of the model increase, define

---

[11] At the cost of lengthening some of the proofs, we could additionally require that $a$ and $a'$ are identical except that at some histories $a'$ initiates a new field while $a$ does not.

$(a_1, ..., a_t, ...)$ to be an *equilibrium sequence* if, for $\tau \geq 1$, the plan of length $\tau$, $(a_1, ..., a_\tau)$, forms an equilibrium (for $T = \tau$). We will say that *with probability* 1 *a sequence underinvests in new fields at some date* if there is a set of states $A$ with $P(A) = 1$ where, for any $\omega \in A$, there is a $L$ such that any plan $(a_1, ..., a_\tau)$ in the sequence with $\tau \geq L$ underinvests in new fields at some date given $\omega$.

# Appendix B: remaining proofs

**Proof of cutoff lemma.** I. Suppose $\beta \geq 2$ and $t \leq T - 2$. To calculate the expected value of research achieved by the equilibrium plan $a$ from $t$ onwards, we first calculate the expected values of research in periods $t$ through $t + 2$.

Let $V_{t+2}(v_{-1})$ be the expected value of research at date $t+2$ in equilibrium given a history $h_{t+1}$ where $v_{-1}$ is the maximum of 0 and the value of the highest-value leading project of $h_{t+1}$. Let $V_{t+1}(v_{-1}^l, v_{-1}^h)$ be the sum of the expected value of research at dates $t + 1$ and $t + 2$ in equilibrium given a $h_t$ where $v_{-1}^h$ (resp. $v_{-1}^l$) is the maximum of the value of the highest value (resp. second-highest value) leading project of $h_t$ and 0, and let $V_t(v_{-1}^l, v_{-1}^m, v_{-1}^h)$ be the sum of the expected value of research from $t$ through $t + 2$ in equilibrium given a $h_{t-1}$ where $v_{-1}^h$ (resp. $v_{-1}^m$, $v_{-1}^l$) is the maximum of the value of the highest value (resp. second-highest value, third-highest value) leading project of $h_{t-1}$ and 0. We have

$$V_{t+2}(v_{-1}) = \delta v_{-1},$$

$$V_{t+1}(v_{-1}^l, v_{-1}^h) = \delta v_{-1}^h + \int_{v_{-1}^l}^{\infty} V_{t+2}(v) g_{v_{-1}^h}(v) dv + \int_{-\infty}^{v_{-1}^l} V_{t+2}(v_{-1}^l) g_{v_{-1}^h}(v) dv,$$

$$V_t(v_{-1}^l, v_{-1}^m, v_{-1}^h) = \delta v_{-1}^h + \int_{v_{-1}^m}^{\infty} V_{t+1}(v_{-1}^m, v) g_{v_{-1}^h}(v) dv$$

$$+ \int_{v_{-1}^l}^{v_{-1}^m} V_{t+1}(v, v_{-1}^m) g_{v_{-1}^h}(v) dv + \int_{-\infty}^{v_{-1}^l} V_{t+1}(v_{-1}^l, v_{-1}^m) g_{v_{-1}^h}(v) dv.$$

These formulas are mostly self-explanatory. The first term in the expression for each $V_\tau$ is the expected value of an immediate successor of the highest-value leading project of $h_{\tau-1}$, which will be project undertaken at $\tau$, while the remaining terms are the expectation of $V_{\tau+1}$ using the realization of the project undertaken at $\tau$ to determine the highest-value leading project of $h_{\tau+1}$. In the expression for $V_{t+1}(v_{-1}^l, v_{-1}^h)$, for example, the second integral indicates the fact that if the successor of the project with value $v_{-1}^h$ turns out to have value less than $v_{-1}^l$ then at date $t + 2$ a successor of the project with value $v_{-1}^l$ will be selected.

Next consider an alternative plan that (i) at $t$, initiates a new field $f'$, (ii) at $t+1$, selects a successor of $(0, f')$ if $v(0, f') > 0$ and otherwise initiates a new field $f''$, and (iii) at $t + 2$, selects a successor of $(0, f')$ or $(0, f'')$ if either $v(0, f') > 0$ or $v(0, f'') > 0$ and otherwise initiates a new field. Notice that with this strategy the expected values of research at $t$, $t + 1$, and $t + 2$ are not functions of the values of the leading projects of $h_t$ and the expected values at $t + 1$ and $t + 2$ are functions only of $v(0, f')$. So, letting $\widehat{V}_\tau$ denote the sum of the expected values of research from $t$ through $t + 2$ for this alternative strategy and, at both $t + 2$ and $t + 1$, letting $v_{-1}$ denote the maximum expected value of leading projects of $h_\tau$ in

fields $f'$ or $f''$, we have

$$\widehat{V}_{t+2}(v_{-1}) = \delta v_{-1},$$

$$\widehat{V}_{t+1}(v_{-1}) = \delta v_{-1} + \int_{v_{-1}}^{\infty} \widehat{V}_{t+2}(v)g_{v_{-1}}(v)dv + \int_{-\infty}^{v_{-1}} \widehat{V}_{t+2}(v_{-1})g_{v_{-1}}(v)dv,$$

$$\widehat{V}_t = \int_0^{\infty} \widehat{V}_{t+1}(v)g(v)dv + \int_{-\infty}^{0} \widehat{V}_{t+1}(0)g(v)dv.$$

The key item above is the second integrand in the expression for $\widehat{V}_{t+1}(v_{-1})$, which indicates that if the root project in the field initiated at $t$ has value $v_{-1}$ and the successor project undertaken at $t+1$ turns out to have value less than $v_{-1}$ then at $t+2$ a project in a different branch of the new field will be undertaken with an expected value of $\widehat{V}_{t+2}(v_{-1})$.

We now compare the sum of expected values of research from $t$ through $t+2$ for these two strategies – $V_t(v_{-1}^l, v_{-1}^m, v_{-1}^h)$ versus $\widehat{V}_t$ – given various candidate cutoffs that we label as $\varepsilon(i)$. Let $\langle \varepsilon(n) \rangle$ be a sequence of strictly positive numbers such that $\varepsilon(n) \to 0$ and let $\langle v_{-1}^l(n), v_{-1}^m(n), v_{-1}^h(n) \rangle$ be a sequence of triples such that, for all $n$, $0 \leq v_{-1}^k \leq \varepsilon(n)$ for $k = l, m, h$ and $v_{-1}^l(n) \leq v_{-1}^m(n) \leq v_{-1}^h(n)$, indicating the three highest values of leading projects at $t$, except that as before, since negative expected-value projects are not undertaken, 0's replace negative values. Then $\delta v_{-1}^h(n) \to 0$ and $\int_{v_{-1}^l(n)}^{v_{-1}^m(n)} V_{t+1}(v, v_{-1}^m(n))g_{v_{-1}^h(n)}(v)dv \to 0$ as $n \to \infty$. Substituting in the definitions of $V_{t+1}(v_{-1}^l, v_{-1}^h)$, $V_{t+2}(v_{-1})$, $\widehat{V}_{t+1}(v_{-1})$, $\widehat{V}_{t+2}(v_{-1})$, it is readily confirmed that

$$\int_{-\infty}^{v_{-1}^l(n)} V_{t+1}(v_{-1}^l(n), v_{-1}^m(n))g_{v_{-1}^h(n)}(v)dv \to \int_{-\infty}^{0} \widehat{V}_{t+1}(0)g(v)dv.$$

Consider finally the remaining term in $V_t(v_{-1}^l(n), v_{-1}^m(n), v_{-1}^h(n))$,

$$\int_{v_{-1}^m(n)}^{\infty} V_{t+1}(v_{-1}^m(n), v)g_{v_{-1}^h(n)}(v)dv = \int_{v_{-1}^m(n)}^{\infty} \left[ \delta v + \int_{v_{-1}^m(n)}^{\infty} \delta \widetilde{v} g_v(\widetilde{v})d\widetilde{v} \right. \tag{1}$$

$$\left. + \int_{-\infty}^{v_{-1}^m(n)} \delta v_{-1}^m(n)g_v(\widetilde{v})d\widetilde{v} \right] g_{v_{-1}^h(n)}(v)dv,$$

and compare it to the remaining term in $\widehat{V}_t$,

$$\int_0^{\infty} \widehat{V}_{t+1}(v)g(v)dv = \int_0^{\infty} \left[ \delta v + \int_v^{\infty} \delta \widetilde{v} g_v(\widetilde{v})d\widetilde{v} + \int_{-\infty}^{v} \delta v g_v(\widetilde{v})d\widetilde{v} \right] g(v)dv \tag{2}$$

$$= \int_0^{\infty} \left[ \delta v + \int_v^{\infty} \delta \widetilde{v} g_v(\widetilde{v})d\widetilde{v} + \int_0^{v} \delta v g_v(\widetilde{v})d\widetilde{v} + \int_{-\infty}^{0} \delta v g_v(\widetilde{v})d\widetilde{v} \right] g(v)dv.$$

Now, for any $v \geq 0$ and any $v_{-1}^m(n) \geq 0$, $\int_v^{\infty} \delta \widetilde{v} g_v(\widetilde{v})d\widetilde{v} + \int_0^{v} \delta v g_v(\widetilde{v})d\widetilde{v} > \int_{v_{-1}^m(n)}^{\infty} \delta \widetilde{v} g_v(\widetilde{v})d\widetilde{v}$,

32

which implies

$$\lim_{n\to\infty}\int_0^\infty\left[\int_v^\infty \delta\tilde{v}g_v(\tilde{v})d\tilde{v}+\int_0^v \delta vg_v(\tilde{v})d\tilde{v}\right]g(v)dv\geq\lim_{n\to\infty}\int_{v_{-1}^m(n)}^\infty\left[\int_{v_{-1}^m(n)}^\infty \delta\tilde{v}g_v(\tilde{v})d\tilde{v}\right]g_{v_{-1}^h(n)}(v)dv.$$

Since in addition

$$\int_{v_{-1}^m(n)}^\infty \delta vg_{v_{-1}^h(n)}(v)dv\;\;\to\;\;\int_0^\infty \delta vg(v)dv,\text{ and}$$

$$\int_{v_{-1}^m(n)}^\infty\left[\int_{-\infty}^{v_{-1}^m(n)} \delta v_{-1}^m(n)g_v(\tilde{v})d\tilde{v}\right]g_{v_{-1}^h(n)}(v)dv\;\;\to\;\;0,$$

the difference between (2) and (1) converges to a number at least as great as $\int_0^\infty[\int_{-\infty}^0 \delta vg_v(\tilde{v})d\tilde{v}]g(v)dv$, a strictly positive constant. Hence for all $n$ sufficiently large and hence all $\varepsilon(n)$ sufficiently small $\widehat{V}_t > V_t(v_{-1}^l(n), v_{-1}^m(n), v_{-1}^h(n))$, which shows that for all sufficiently large $n$ the alternative strategy delivers larger expected value of research from $t$ to $t+2$ when $h_{t-1}$ is such that the highest value leading project has a value in $(0, \varepsilon(n))$.

We fill in the remainder of the alternative plan by defining project choices for periods $t+3, ..., T$ that will, as $\varepsilon(n) \to 0$, yield an expected value of research in these periods that converges to the expected value achieved by the equilibrium plan $a$ in the same periods. For $k = 1, ..., T$, let $z_k$ denote a realization of the deviation of $v(\pi)$ from its expected value where $\pi$ is the project undertaken in equilibrium at $k$ when the deviations $z_1, ..., z_{k-1}$ have been realized. Given a history $h_{t+2}$, which is uniquely defined by the equilibrium plan $a$ and the realizations $z_1, ..., z_{t+2}$, suppose the equilibrium has $r$ projects feasible at $t+3$ with strictly positive expected value. Each of these projects must be in a distinct branch, which we label $1, ..., r$, of a single field $\widehat{f}$. It will be convenient to henceforth label the projects in these branches so that $(j, i, \widehat{f})$ for $i \geq 0$ now denotes the $(i+1)$th successor in branch $j$ of the leading project at $t+3$ of the $j$th branch of $\widehat{f}$.[12] Let $f_1, ..., f_r$ index $r$ fields that the equilibrium plan $a$ has not initiated by date $t+2$. We may then define a 'preliminary' alternative plan $\alpha'$ that undertakes project $(1, i, f_j)$ at period $t'$, given that the projects undertaken earlier by $\alpha'$ have realized deviations $z_1, ..., z_{t'-1}$, whenever the equilibrium $a$, given the same deviations $z_1, ..., z_{t'-1}$, undertakes $(j, i, \widehat{f})$. Finally, let $\alpha'$, given the realized deviations $z_1, ..., z_{t'-1}$, initiate a new field $f_{alt}$ (where $f_{alt} \notin \{f_1, ..., f_r\}$) whenever the equilibrium following the same deviations initiates a new field $f_{eq}$, and then, for $t'' > t'$, given the deviations $z_1, ..., z_{t''-1}$ undertake $(b, i, f_{alt})$ whenever the equilibrium given the same deviations undertakes $(b, i, f_{eq})$. Define a project undertaken at date $i$ to have a 'hypothetical value' equal to $\delta v + z_i$ when its predecessor has value $v$, even when $v < 0$. We will now see that, given any $z = (z_1, ..., z_T)$, the sum of hypothetical values that $\alpha'$ generates from $t+3$ to $T$ converges, as $\varepsilon(n) \to 0$, to the expected value generated by the equilibrium from $t+3$ to $T$. It follows (see part II below) that the plan $\alpha$ that is identical to $\alpha'$ except that $\alpha$ initiates a new field whenever $\alpha'$ undertakes a nonpositive expected value

---

[12]So for example $(j, 0, \widehat{f})$ is the immediate successor in branch $j$ of the branch $j$ leading project of $\widehat{f}$. Note that the leading project of branch $j$ may be $(0, \widehat{f})$.

project must then both generate greater total expected value than the equilibrium $a$.

To conclude part I, therefore, we show that as $\varepsilon(n) \to 0$ the expected value of research from $t + 3$ through $T$ under the equilibrium $a$ converges to the sum of the hypothetical expected values for the same periods under $\alpha'$. Since fixed dispersion implies that each measurable set of deviation vectors $z = (z_1, ..., z_T)$ has the same probability in the two strategies, it is sufficient to show that, for any $z$, the difference between the expected value of research from $t + 3$ to $T$ delivered by $a$ and the hypothetical values delivered by $\alpha'$ in the same periods is bounded above by $\varepsilon(n)(1 + \delta + ... + \delta^{T-(t+3)}) = \varepsilon(n) \sum_{l=0}^{T-(t+3)} \delta^l$. The first $t + 2$ coordinates of $z$ determine $h_{t+2}$ which we may now take as fixed.

Let $e$ denote $E(v(j, 0, \widehat{f})|h_{t+2})$ and let $\eta = (\eta_0, \eta_1, ...)$ be a sequence of deviations of project values from their expected values for projects $((j, 0, \widehat{f}), (j, 1, \widehat{f}), ...)$. Define $(v(j, i, \widehat{f}))_{i \geq 0}$ recursively by $v(j, 0, \widehat{f}) = e + \eta_0$ and $v(j, k, \widehat{f}) = \delta v(b, k - 1, \widehat{f}) + \eta_k$ for $k \geq 1$. Also, given $\eta$ and $e$, define $S_k(\eta, e) = \sum_{l=0}^k v(j, l, \widehat{f})$. It is easy to see that $S_k(\eta, e) = e(1 + \delta + ... + \delta^k) + \eta_1(1 + \delta + ... + \delta^k) + \eta_2(1 + \delta + ... + \delta^{k-1}) + ... + \eta_k$.

As for the preliminary alternative, if $\alpha'$ initiates $f_i$ at some period and $\eta_k$ is a realization of the deviation of $\nu(1, k, f_i)$ from its expected value then, given the realizations $\eta = (\eta_0, \eta_1, ...)$, define the hypothetical values $(w(1, i, f_i))_{i \geq 0}$ recursively by $w(1, 0, f_i) = \eta_0$ and $w(1, k, f_i) = \delta w(1, k - 1, f_i) + \eta_k$ for $k \geq 1$. For any $\eta$, $\sum_{l=0}^k w(1, l, f_i) = S_k(\eta, 0)$. So, for any $\eta$ and $k$, $S_k(\eta, e) - S_k(\eta, 0)$ equals $e \sum_{l=0}^k \delta^l$.

Now for an arbitrary $z = (z_1, ..., z_T)$ and given the equilibrium plan $a$, some subset of the coordinates of $z$ will be the deviations $(\eta_0, ..., \eta_\tau)$ for the projects in branch $j$ of $\widehat{f}$. Since $\alpha'$ undertakes projects in $f_j$ if and only if $a$ undertakes projects in branch $j$ of $\widehat{f}$, the deviations for the projects undertaken by $\alpha'$ in $f_j$ will be $(\eta_0, ..., \eta_\tau)$ when $z$ obtains. Hence given $z$ and $\varepsilon(n)$, the difference between the values delivered by $a$ from $t + 3$ to $T$ and the hypothetical values delivered by $\alpha'$ is indeed bounded above by $\varepsilon(n) \sum_{l=0}^k \delta^l$.

II. To conclude we assume that $\beta = 1$ or $t \in \{T - 1, T\}$ and show that initiating a new field rather undertaking a feasible project in an existing field with positive expected value cannot increase the total expected value of research. Observe first that with $\beta = 1$ or $t \in \{T - 1, T\}$ it is impossible to initiate a field $f$ at $t$ and then undertake projects in more than one distinct branch of $f$. Now consider plans that begin by selecting a project with expected value $v$, that always select from only one branch of any given field, and that initiate a new field immediately following the selection of a negative value project. It is easy to confirm that if $v > 0$ then for any integer $n > 0$ the distribution of values for the $n$th selection of such a plan will first order stochastically dominate a plan that begins a project with $v = 0$. Given that the only feasible choices in any $f$ must be drawn from a single branch, total expected value is maximized by selecting the project with the highest expected value. ∎

**Proof of Proposition 5.** Let $a$ be an optimal plan, which achieves the total expected value $EV$, and let $n$ be the number of fields whose root projects are undertaken at some history with plan $a$. No projects are bundled at history $\varnothing$ and set $p_\varnothing(0, f) = \frac{EV}{n}$ if $f$ is undertaken at some history and $p_\varnothing(0, f) = 0$ otherwise. Agent 1 buys each $(0, f)$ and undertakes $a(\varnothing)$. If $h_1$ obtains, agent 1 bundles all feasible projects into a single bundle, which will sell at a price equal to the conditional total expected value of $a$ given $h_1$. Each

subsequent agent $i$ buys the one bundle created at $i-1$, undertakes $a(h_{i-1})$, again, if $h_i$ obtains, bundles all feasible projects which sell for the conditional total expected value of $a$ given $h_i$. All other bundles sell for price 0. ∎

**Proof of Proposition 6.** Given a value-seeking equilibrium plan $a$, suppose there is a last date $t$ and accompanying history $h_{t-1}$ at which a citation-maximizer could increase his expected citations by undertaking a project $\pi' = (1, i', f')$ that differs from $\pi = a(h_{t-1}) = (1, i, f)$. Let $a'$ be the plan that coincides with $a$ except that $a'(h_{t-1}) = \pi'$. Since in a value-seeking equilibrium at most one project can have a strictly positive conditional expected value given $h_{t-1}$, and if there is such a project it must be undertaken, $E(v(\pi')|h_{t-1}) = 0$ and $E(v(\pi)|h_{t-1}) \geq 0$. Let $\pi'_j$ denote $(1, i'+j, f')$ and $\pi'_j$ denote $(1, i+j, f)$ for $j = 0, ..., n$. Given the random walk assumption, any particular sequence $(v(\pi'_1), ..., v(\pi'_n))$ of possible realizations can be viewed as a sequence $s = (s_1, ..., s_n)$ of successful and unsuccessful outcomes of an unbiased coin: $v(\pi'_k)$ is a success and we set $s_k = 1$ (resp. failure and $s_k = -1$) if and only if $v(\pi'_k) - v(\pi'_{k-1}) > 0$ (resp. $< 0$). Each sequence $s$ has probability $\frac{1}{2^n}$. If $T - t \geq n$, then for agent $t$ to earn exactly $n$ citations from undertaking $\pi'$, we must have $v(\pi'_j) \geq 0$ for $j = 1, ..., n-1$, since if $v(\pi'_j) < 0$ no value-maximizer would undertake $\pi'_{j+1}$. Equivalently the cumulative total of successes $\sum_{k=1}^{m} s_k$ must not drop below 0 for $m = 1, ..., n-1$. Let $S$ be the set of sequences defined by $s = (s_1, ..., s_n) \in S$ if and only if (1) $n \leq T - t$, (2) $\sum_{k=1}^{m} s_k \geq 0$ for $m = 1, ..., n-1$, and (3) if $n < T - t$ then $\sum_{k=1}^{n} s_k = 0$.

Given that behavior after $t$ is value-seeking, the conditional probability that $t$ earns exactly $n$ citations with plan $a'$, given that $(v(\pi'_1), ..., v(\pi'_n))$ is determined by the sequence $s \in S$, may be less than 1. If on the other hand $t$ undertakes $\pi$ and plan $a$ obtains, and $(v(\pi_1), ..., v(\pi_n))$ is determined by the same $s$, then $v(\pi_j) > 0$ for $j = 1, ..., n-1$ and therefore the conditional probability that $t$ earns at least $n$ citations given $s$ equals exactly 1.

Therefore, letting $P_{a'}(n, s)$ be the conditional probability that $t$ earns exactly $n$ citations under plan $a'$, given that $(v(\pi'_1), ..., v(\pi'_n))$ is determined by $s \in S$, letting $P_a(n, s)$ be the conditional probability that $t$ earns at least $n$ citations under plan $a$, given that $(v(\pi_1), ..., v(\pi_n))$ is determined by $s \in S$, and letting $n(s)$ be the number of entries in $s$, we have $P_{a'}(n(s), s)n(s) \leq P_a(n(s), s)n(s)$ for any $s \in S$. Hence

$$\sum_{s \in S} \frac{1}{2^{n(s)}} P_{a'}(n(s), s)n(s) \leq \sum_{s \in S} \frac{1}{2^{n(s)}} P_a(n(s), s)n(s).$$

Since the number of citations earned by $t$ given plan $a$ and $s \in S$ is greater than or equal to $n(s)$, the expected number of citations with $\pi$ must be at least as great as with $\pi'$.

Hence the sum from $n = 1$ to $n = T$ of $n$ times the probability of the sequences of successes and failures at which $\pi'$ earns exactly $n$ citations can be no larger than the sum from $n = 1$ to $n = T$ of $n$ times the probability of the sequences of successes and failures at which $\pi$ earns exactly $n$ citations for $n = 1, ..., T$. ∎

**Proof of Proposition 7.** Consider the citations equilibrium $a$ and the positive-probability set of histories $H'$ through period $T - 5$ such that, for all $t \leq T - 5$ and all subhistories $h_{t-1}$ of $h \in H'$, $v(a(h_{t-1})) < 0$. Let us assume that $a(h_{T-5})$ is a root project for a positive-probability $H'' \subset H'$, since otherwise the Proposition is proved. We set the set of histories $H$ in the Proposition to consist of all histories through $T - 4$ that continue some $h \in H''$ and

35

where $v(a(h_{T-5})) \in (0, \varepsilon)$ for all $h_{T-5} \in H''$, for a $\varepsilon > 0$ to be specified later, and assume that at equilibrium $a$ and all continuations of $h \in H$ agents $T - 2$ through $T$ maximize value since again otherwise the Proposition would be proved.

It is sufficient to show that, for one of the histories $h_{T-4}$ we have identified, $a(h_{T-4})$ is not a successor of $a(h_{t-5})$. Suppose $T - 3$ undertakes a root project $\pi$ following $h_{T-4}$. If $v(\pi) > \varepsilon$ then, since $\beta \geq 4$ and agents $T - 2$ through $T$ maximize value, $T - 2$ through $T$ will undertake successors of $\pi$. Since $P(v(\pi) > \varepsilon) = \frac{1}{2} \int_{\varepsilon}^{\infty} g(v) dv$ converges to $\frac{1}{2}$ as $\varepsilon$ converges to 0, agent $T - 3$ will earn nearly $\frac{3}{2}$ expected citations with $\pi$ for small enough $\varepsilon$. If $T - 3$ undertakes a successor $\pi'$ of $a(h_{T-5})$ and if, using the terminology of the proof of Proposition 6, $\pi'$ and its next two successors are all successes, which we write at SSS, then $T - 3$ earns 3 citations. If SSF then $T - 3$ earns 3 citations, if SFS then in some citations equilibria $T - 3$ can earn as many as 3 citations, if SFF then $T - 3$ can earn as many as 2 citations at some equilibria, and finally if $\pi'$ is a failure, then, since the remaining agents maximize value and $\beta \geq 4$, $T - 3$ earns 0 citations. Since each combination of three outcomes has probability $\frac{1}{8}$, $T - 3$ by choosing a non-root project can earn at most $\frac{1}{8}(3 + 3 + 3 + 2) < \frac{3}{2}$ expected citations. So, if we set $\varepsilon > 0$ so that $\frac{3}{2} \int_{\varepsilon}^{\infty} g(v) dv > \frac{11}{8}$ then $T - 3$ will earn more expected citations by undertaking $\pi$ than by undertaking any successor of $a(h_{T-5})$. ∎

**Proof of Proposition 8.** Following the proof of Proposition 7, consider a plausible citations equilibrium $a$ and the positive-probability set of histories $H'$ through period $T - 5$ such that for all subhistories $h_t$ of $h \in H'$ with $t < T - 5$, $v(a(h_t)) < 0$. Due to plausibility, no agent $t \in \{2, ..., T - 5\}$ at any subhistory $h_{t-1}$ of any $h \in H'$ chooses a successor of any $a(h_\tau)$, $\tau < t$. To see why, observe that if an agent were to undertake the successor $\pi'$ of a negative value project then by plausibility all subsequent agents who undertake successors of $\pi'$ must undertake negative expected value projects, and the latest date at which which an agent undertakes such a successor at some history would gain strictly more expected citations by undertaking a project with the highest available expected value. So, for all $h_{T-5} \in H'$, $a(h_{T-5})$ is a root project. We restrict the set of histories $H$ to the continuations of $h_{T-5} \in H'$ such that $v(a(h_{T-5})) \in (0, \varepsilon)$, for a $\varepsilon > 0$ such that $\frac{3}{2} \int_{2\varepsilon}^{\infty} g(v) dv > \frac{11}{8}$.

Given Proposition 1, for sufficiently small $\varepsilon > 0$, the initiation of a new field at $h_{T-4} \in H$ will increase the conditional total expected value produced from $T - 3$ through $T$ given $h_{T-4}$. And given the proof of Proposition 7, it remains to show only that if $T - 3$ initiates a new field $\pi$ with $v(\pi) > 2\varepsilon$ then $T - 2$ and $T - 1$ will earn strictly more expected citations from choosing successors of $\pi$ than from any other feasible choice. So assume henceforth that $T - 3$ initiates a new field $\pi$ with $v(\pi) > 2\varepsilon$.

First, since $T$ maximizes value and if tied value are 0 probability events, $T - 1$ must maximize the likelihood that $v(a(h_{T-2}))$ has a value that is strictly greater than any other available project. Next, to see that $T - 2$ will choose a successor $\pi_{T-2}$ of $\pi$, observe that if $T - 2$ does so and $\pi_{T-2}$ is a success (i.e., $v(\pi_{T-2}) > v(a(\pi))$) then $T - 1$ will maximize the likelihood that $v(a(h_{T-2}))$ is the highest-value project by choosing the successor $\pi_{T-1}$ of $\pi_{T-2}$ (given that $v(\pi) > 2\varepsilon$ and given that the random walk assumption implies that the probability that the value of a new field will be greater than or equal to $v(\pi_{T-2})$ is less than $\frac{1}{2}$). Similarly, if $\pi_{T-2}$ and $\pi_{T-1}$ are both successes then $T$ must choose a successor of $\pi_{T-1}$. So if $T - 2$ chooses a successor of $\pi$ then he earns at least $\frac{3}{4}$ expected citations. If in contrast $T - 2$ chooses any other project $\pi'$ then $P(v(\pi') > 2\varepsilon) < \frac{1}{2}$ and $P\left(v(\pi') \geq v(\pi) | v(\pi) > 2\varepsilon \text{ and } v(\pi') > 2\varepsilon\right) = \frac{1}{2}$,

the latter due to symmetry. Hence $P(v(\pi') \geq v(\pi)|v(\pi) > 2\varepsilon) < \frac{1}{4}$. When $v(\pi') \geq v(\pi)$, $T - 2$ earns at most 2 citations and when $v(\pi') < v(\pi)$, $T - 2$ earns 0 citations since $\beta \geq 3$ and $T - 1$ and $T$ are value-maximizers. Hence $T - 2$ earns less than $\frac{1}{2}$ expected citations and will therefore undertake a successor $\pi_{T-2}$ of $\pi$. Given this decision for $T - 2$, $T - 1$ will also choose a successor $\pi_{T-2}$ of $\pi$. ∎

**Proof of Proposition 9.** Let $a$ be a plausible citations equilibrium. Given $\varepsilon > 0$, let $\Delta = \{\ \delta \in (1 - \varepsilon, 1) : \int_0^\delta g_n(v)dv < \varepsilon\}$ which is nonempty for all $n$ sufficiently large.

Let the positive-probability set $H$ be a set of histories of length $T - 4$ such that (1) $v(a(h_t)) < 0$ for all subhistories $h_t$ of $h \in H$ that have length $0 \leq t \leq T - \beta - 5$, (2) $v(a(h_{T-\beta-4})) \in \Delta$ for any subhistory $h_{T-\beta-4}$ of $h \in H$ of length $T - \beta - 4$, and where by plausibility $a(h_{T-\beta-4})$ must be a root project, say $(0, f)$, and (3) $v(b, 1, f) = 0$ for any $b \in \{1, ..., \beta\}$ such that $(b, 1, f) = a(h_t)$ at some subhistory $h_t$ of $h \in H$ such that $T - \beta - 3 \leq t \leq T - 5$.

There are two types of plausible equilibria at histories in $H$, those where $a(h_t)$ is in $f$ for all subhistories $h_t$ of $h \in H$ with length $T - \beta - 3 \leq t \leq T - 5$ and those where $a(h_t)$ is a root project at one or more of the named subhistories. To see that at equilibria of the first type, $a(h_{T-4})$ is a root project for all $h_{T-4} \in H$, observe first that if $\pi$ is a root project at $h_{T-4}$ then as $n$ increases $P(v(\pi) > v(a(h_{T-\beta-4})))$ converges to $\frac{1}{2}$. Now if $a(h_{T-4})$ is a root project $\pi$ and $P(v(\pi) > v(a(h_{T-\beta-4})))$ then $T - 3$ earns 3 citations at $h_{T-4}$. To see this, observe that if $T - 2$ initiates a new field $f'$ then, with probability less than $\frac{1}{2}$, $v(0, f')$ will have a value high enough for agent $T$ to choose a project in $f'$ even if $T - 1$ chooses a project in $f'$, giving $T - 2$ less than 1 expected citations, while if $T - 2$ chooses a successor of $\pi$, say $\pi'$, and $\pi'$ is a success then $T - 2$ earns 2 citations while if $\pi'$ is a failure then $T - 2$ earns 0 citations (due to plausibility), giving $T - 2$ exactly 1 expected citation. Hence if $T - 3$ chooses a root project he earns nearly $\frac{3}{2}$ expected citations when $n$ is large. If on the other hand $a(h_{T-4})$ is a successor $(b, 1, f)$ of $a(h_{T-\beta-4})$, then $P(v(b, 1, f) = 0|h_{T-4}) = \frac{1}{2}$ and $P(v(b, 1, f) > 0, v(b, 2, f) > 0, \text{ and } v(b, 1, f) = 0|h_{T-4}) = \frac{1}{8}$. Given plausibility, it follows that if $a(h_{T-4})$ is a successor of $a(h_{T-\beta-4})$ then $T - 3$ earns strictly less than $\frac{3}{2}$ citations at $h_{T-4}$ (for all $n$).

Thus at both types of plausible equilibria, citation maximizers at some history or sub-history $h$ in $H$ of length $t$ initiate a root project rather than undertake a project with an expected value in $\Delta$. We adopt the notation used in the proof of Proposition 6. Given any sequence $s = (s_1, ..., s_{T-t})$ of $T - t$ successes and failures of the projects undertaken following $h$, any value-maximizing plan $a_v$ produces a conditional expected total value after $t$ given $s$ that converges, as $n \to \infty$, to the conditional expected total value after $t$ given $s$ produced by $a$. To see why, we can put aside the difference between the value increment of a new field and that of $f$, which converges to 0 as $n \to \infty$. Since for any $(b, 1, f)$ not undertaken through $h$ satisfies $E(v(b, 1, f)|h) > 0$ and by assumption $E(v(a(h))|h) = 0$, the value realized given $s$ by $a_v$ is strictly greater than the value realized by $a$ up until and including the first date $T - 4 + m_1$ such that $\sum_{k=1}^{m_1} s_k \leq 0$. If at the first $m_i$ such that $\sum_{k=m_{i-1}+1}^{m_i} s_k \leq 0$ and both plans have undertaken all of the successors of $(0, f)$, and assuming that $a$ is value-maximizing after $t + 1$, plan $a$ will increase its value produced by 1 relative to $a_v$ but it cannot overtake $a$. For some sequences, however, specifically those with $\sum_{k=1}^m s_k > 1$ for all $m$, $a_v$ produces a conditional expected total value after $t$ given $s$ that, for all $n$ sufficiently large, is discretely

37

greater than the conditional expected total value after $t$ given $s$ produced by $a$. Since the probabilities of each sequence, $\frac{1}{2^{T-t}}$, is strictly positive and equal for both plans, $a_v$ produces greater conditional expected value given $h$ than $a$. ∎

# References

[1] Banerjee, A., 1992, 'A simple model of herd behavior,' *Quarterly Journal of Economics* 107: 797-818.

[2] Bikhchandani, S., Hirshleifer, D. and Welch, I., 1992, 'A theory of fads, fashion, custom, and change as informational cascades,' *Journal of Political Economy* 100: 992-1026.

[3] Dasgupta, P. and David, P., 1994, 'Toward a new economics of science,' *Research Policy* 23: 487-521.

[4] Dasgupta, P. and Maskin, E., 1987, 'The simple economics of research portfolios,' *The Economic Journal* 97: 581-595.

[5] Jovanovic, B. and Nyarko, Y., 1996, 'Learning by doing and the choice of technology,' *Econometrica* 64: 1299-1310.

[6] Hoddeson, L., Brown, B., Riordan, M., and Dresden, M. (eds.), 1997, *The Rise of the Standard Model: A History of Particle Physics from 1964 to 1979*, Cambridge University Press, Cambridge.

[7] Hong, L., and Page, S., 2004, 'Groups of diverse problem solvers can outperform groups of high-ability problem solvers,' *Proceedings of the National Academy of Sciences* 101: 16385-16389.

[8] Kuhn, T., 1962, *The Structure of Scientific Revolutions*, University of Chicago Press: Chicago.

[9] Kuznets, S., 1930, *Secular Movements in Production and Prices*, Houghton-Mifflin, Boston.

[10] McCall, J., 1970, 'Economics of information and job search,' *Quarterly Journal of Economics* 84: 113-126.

[11] Ne'eman, Y., 1982, 'Gauge-theory ghosts and ghost-gauge theories,' in *Differential Geometric Methods in Mathematical Physics, Clausthal 1980* (H.-D. Doebner, S. Andersson, and H. Petry, eds.), Lecture notes in mathematics 905, Springer, Berlin, p. 241-259.

[12] Polyakov, A., 1997, 'A view from the island,' in *The Rise of the Standard Model: A History of Particle Physics from 1964 to 1979* (Hoddeson et al., eds.), Cambridge University Press, Cambridge, p. 243-249.

[13] Rothschild, M., and Stiglitz, J., 1970, 'Increasing risk: I. a definition,' *Journal of Economic Theory* 2: 225-243.

[14] Schumpeter, J., 1911, *The Theory of Economic Development*, Harvard University Press, Cambridge MA, 1934.

[15] 't Hooft, G., 1999, 'When was asymptotic freedom discovered? or The rehabilitation of quantum field theory,' *Nuclear Physics B* 74: 413-425.

[16] Williams, D., 1991, *Probability with Martingales*, Cambridge University Press, Cambridge.