

# THE DOUBLE-EDGED SWORD OF RECOMBINATION IN BREAKTHROUGH INNOVATION

Sarah Kaplan\*  
University of Toronto, Rotman School  
105 St. George Street  
Toronto, ON, M5S 3E6, Canada  
416-978-7403  
[sarah.kaplan@rotman.utoronto.ca](mailto:sarah.kaplan@rotman.utoronto.ca)

Keyvan Vakili  
London Business School  
Sussex Place, Regent's Park  
London, NW1 4SA, UK  
079 6690 8742  
[kvakili@london.edu](mailto:kvakili@london.edu)

This draft: February 5, 2014

*\*Corresponding author. We are grateful to Michael Lounsbury for the suggestion to study fullerenes and nanotubes within the broad domain of nanotechnology, Juan Alcácer for his help in developing the initial data set, and Hanna Wallach for her early test-driving of topic modeling on these data. Grace Gui, Illan Kramer, Sara Sojung Lee, Octavio Martinez, Shamsi Mohtashim, Neal Parikh, Navid Soheilnia, Aaron Sutton, and Xihua Wang provided much-appreciated research assistance. Earlier drafts of this manuscript benefited from the comments of Ajay Agrawal, Christian Catalini, Kristina Dahlin, Andrea Fosfuri, Alfonso Gambardella, Constance Helfat, Nico Lacetera, Anita McGahan, Andrea Mina, Jasjit Singh, Kevyn Yong, and participants in the DRUID conference 2012 (where it received the Best Paper Prize), EGOS Colloquium 2011, the West Coast Research Symposium 2011 and the Rotman Strategy Brownbag. This work is partially supported by the Mack Institute for Innovation Management at the Wharton School, University of Pennsylvania (where Kaplan is a Senior Fellow) and the Canadian Social Sciences and Humanities Research Council under grant #410-2010-0219. All errors and omissions remain our own.*

## THE DOUBLE-EDGED SWORD OF RECOMBINATION IN BREAKTHROUGH INNOVATION

### **Abstract:**

We explore the double-edged sword of recombination in generating breakthrough innovation: recombination of distant knowledge is needed because knowledge in a narrow domain might trigger myopia; but, recombination can be counterproductive when local search is needed to identify anomalies. We take into account how creative processes shape both the cognitive novelty of the idea in addition to the subsequent realization of economic value. We develop a text-based measure of novel ideas in patents by identifying those patents that originate new topics in a body of knowledge using a computer science technique called topic modeling. This measure allows us to distinguish inventions that are novel from those that are valuable (as measured by subsequent citations). We find that, counter to theories of recombination, patents that originate new topics are more likely to be associated with local search, while economic value is the product of broader recombinations as well as novelty.

**Keywords:** breakthrough innovation; recombination; patents; creativity; topic modeling; text analysis; nanotechnology; cognition

# THE DOUBLE-EDGED SWORD OF RECOMBINATION IN BREAKTHROUGH INNOVATION

## INTRODUCTION

Research on innovation has long sought to determine the sources of innovative breakthroughs because they are the basis of change in scientific and technological ideas and of potential increases in social (Kuhn 1962/1996), firm (Hall *et al* 2005, Phene *et al* 2006) and individual economic value (Ahuja *et al* 2005). Most of this research has drawn on theories of recombination based in a ‘tension’ view of the relationship between knowledge and creativity (Weisberg, 1999): deep knowledge in one domain dampens creativity by entrenching researchers into one way of thinking. Breakthrough innovation therefore requires a broad search for information and the recombination of different kinds of knowledge to break those bonds and produce novel ideas that achieve high economic value (Ahuja & Lampert 2001; Guilford 1967). Bridging distant or diverse knowledge or providing structures that enable such recombination should therefore enhance creativity (e.g., Hargadon & Sutton 1997; Audia & Goncalo 2007).

A contrasting, and less tested, theory of creativity – the ‘foundational’ view – differs from the tension view in arguing that local search to identify anomalies is most likely to produce breakthrough innovations (Weisberg 1999; Taylor & Greve 2006). That is, in order to break out of existing constraints and advance a field beyond its current state, one must have a deep understanding of the foundations of a particular knowledge domain, its assumptions and its potential weaknesses. Recombination may be detrimental to innovation because only a deep dive can produce breakthroughs. In theory, all innovations are based on some sort of recombination. We follow the lead of other innovation scholars in referring to recombination as that involving distant or diverse knowledge, where recombination of local or similar knowledge should be seen as local search (e.g., Fleming 2001; Ethiraj & Levinthal 2004).

Rather than seeing the tension and foundational views of recombination as alternatives, we might usefully conceptualize them jointly as a “double-edged sword” (Sternberg & O’Hara 1999: 256):

recombination of distant or diverse knowledge is required to create breakthroughs because knowledge in a narrow domain might trigger intellectual lock in and lead to incremental innovation; on the other hand, such recombination might be counterproductive because local search in a domain may be required in order to find the openings for breakthrough innovations. Yet, to date, the field of management has not examined how the two blades of the sword are interrelated. In our study, we develop a new method for examining these relationships in the context of patented innovations.

We can reconcile the tension and foundational views of recombination and understand their interrelationships if we take into account how creative processes shape both the novelty of the idea (the cognitive dimension of breakthroughs) as well as the subsequent realization of value (the economic dimension of breakthroughs). Innovations potentially differ along these two dimensions. (Amabile 1983; Audia & Goncalo 2007; Amabile & Pillemer 2012). An innovation might be novel in the sense that it introduces a potential new technological trajectory or incremental because it follows on an existing trajectory. Subsequently, an innovation may or may not turn out to be valuable in terms of generating economic returns for the owners of the invention. Art must appeal to collectors or museums, books must be sold to publishers and readers, new startups must attract venture capital funding, patents must garner licensing or sales revenues, etc. Inventors will first generate insights with different levels of novelty. Subsequently, if they understand which ideas will have the most economic value and ‘sell’ the idea to the appropriate audience, the novel ideas will be recognized as valuable (Sternberg & O’Hara 1999). By corollary, if they do not recognize which ideas are most valuable or are unable to sell the ideas, the economic value of the invention will not be realized.

Importantly, while innovation studies have suggested a positive link between novelty and the value of innovations (Trajtenberg *et al* 1997; Singh & Fleming 2010; Phene *et al* 1997), creativity scholars have suggested that the skills associated with generating cognitively novel ideas and those for selecting and promoting ideas that have economic value are weakly related at best (Sternberg & O’Hara 1999; Sternberg 1997). Thus, a comparison of the tension (recombination) and foundational (local search) processes of creativity would benefit from considering their separate impacts on novelty

and value.

In studies of innovation, particularly of innovative patents, research has repeatedly found strong evidence for various forms of recombination as the main mechanism producing breakthroughs (e.g., Fleming 2001; Hall *et al* 2001, Rosenkopf & Nerkar 2001; Hall 2002; Gittelman & Kogut 2003). Yet, the measure of breakthroughs to which they have been constrained is one of citation counts to patents, which have been shown to correlate (if noisily, Bessen 2008) with measures of economic value (Griliches 1990), such as inventors' or other experts' estimates of future financial value (Albert *et al* 1991, Harhoff *et al* 1999), patent renewal fee payments (Harhoff *et al* 1999; Hegde & Sampat 2009), filing patents for the same invention in multiple jurisdictions (Lanjouw & Schankerman 2004), and firms' stock market values (Deng *et al* 1999; Hall *et al* 2005). As a result, citations may most appropriately measure the value or usefulness of patents but do not capture their novelty.

To develop a separate measure of cognitive novelty, we draw on Kuhn's (1962/1996) argument that scientific ideas are embedded in vocabularies and therefore shifts in ideas can be detected in shifts in language. If we are interested in understanding the emergence of breakthrough novel ideas, then, we need methods that pay attention to the language that represents the innovations. In this paper, we introduce just such an approach. Borrowing a computer science technique called 'topic modeling' that discovers the latent topics in a collection of documents and identifies which composition of these topics best accounts for each document, we map the formation of new topics in patent data – which can be seen as the emergence of novel ideas –and, further, locate the patents that introduce them. Those patents that originate new topics can be thought of as cognitive breakthroughs. In introducing this new measure of cognitive breakthroughs, we can unpack the relationship between contrasting creative processes – tension vs. foundational – and differing creative outcomes – cognitive novelty vs. economic value.

To develop and validate this approach, we examine the formation of novel ideas in a domain of nanotechnology, that of Buckminsterfullerenes (and the related area of carbon nanotubes). This is a useful setting because fullerenes can be seen as a 'general purpose technology' (Bresnahan &

Trajtenberg, 1995; Helpman, 1998) with potential applications in many areas (Rothaermel & Thursby, 2007). Thus, a wide variety of cognitive and economic breakthroughs should be possible to identify.

In this nanotechnology field, we find that – consistent with tension theories of creativity – various forms of recombination are positively associated with economic value as measured by patent citation rates. However, in support of foundational theories, we find that cognitive novelty (a patent that originates a new topic) is more likely to be associated with local search that is the product of narrower recombinations. At the same time, we show that novel ideas tend to have higher economic value. This suggests an alternative model of innovation, where novelty is one source of economic value produced by innovators who ‘draw on a single domain in a practiced manner’ (Taylor & Greve 2006, p. 727), while the recombination of distant or diverse knowledge directly positively influences the economic value of innovation, though reducing the likelihood of developing cognitively novel breakthroughs. Few patents in our study were both cognitive and economic breakthroughs (less than one percent of our sample), but they appear to have a greater impact on future innovation than any other kind of invention.

A text-based approach to the analysis of patents gives the researcher new traction in understanding breakthroughs and the emergence and evolution of technologies over time. First, we use texts of patents to develop a measure of cognitive novelty and find that novelty contributes to the creation of economic value. At the same time, we highlight the conflicting creative processes that lead directly to novelty and value, the former requiring local search and the latter distant and diverse recombinations. This contrasts with the common view that local search should be associated with exploitation and not exploration. It also draws attention to the organizational design implications for managing the double-edged sword of recombination in innovation, where innovation strategies must deal with the trade-offs and interrelationships between allocation of resources towards the development of deep knowledge in particular domains and the creation of opportunities for recombination.

## **HYPOTHESES: EXPLORING THE DOUBLE-EDGED SWORD**

The creativity literature proposes that the creative process involves the generation of novelty and then the subsequent achievement of economic value through the recognition and promotion of those novel ideas that have the most promise (Sternberg & O’Hara 1999). This can be seen as a process of variation (either “blind” or intentional production of novelty) followed by selection and retention (realization of economic value) (Campbell 1960; Simonton 1999). Creativity research has proposed two different models – the tension and foundational views – of the role of knowledge in these creative processes. The tension view asserts that deep knowledge can lead to myopia such that recombination of distant or diverse knowledge is needed in order to see new ideas. The foundational view suggests that the only way to see potential anomalies that could lead to breakthroughs is through search in a narrower domain, i.e. local search. These are the two blades of the double-edged sword: recombination is either seen as promoting or detracting from innovation.

Figure 1 portrays the two blades of the sword as they relate to each innovative outcome – cognitive novelty and economic value. In the tension view, recombination should be positively associated with novelty (Path A) and economic value (Path C). In the foundational view, local search is more likely to produce novel (Path A’) and economic value (Path C’). In both cases, novelty, once achieved, should also be associated with economic value (Path B).

--Insert Figure 1 about here --

Note that most innovation studies have an implicit model of innovation based in the tension theory of creativity: they argue that recombination generates novel ideas which, in turn, are more likely to be valuable (in these studies, the focus is patented inventions, so economic value is measured as citations as prior art by subsequent patents) (Fleming 2001; Hall *et al* 2001, Trajtenberg *et al* 1997; Gittelman & Kogut 2003; Singh & Fleming 2010). The dynamics are represented in Figure 1 in Paths A and B. To date, however, these studies have mainly looked at the effect of recombination in patents (using a variety of measures) on subsequent citations to those patents. In other words, rather than testing paths A and B separately, they have tested Path C in Figure 1. They find that recombination is

positively associated with economic value but do not analyze directly the intervening step associated with the generation of novelty from those recombination processes.

In using topic modeling to analyze breakthrough patents, we can explore this relationship directly by measuring the presence of novel ideas (as indicated by shifts in vocabularies in the patent texts) and determining if this variable mediates the association between recombination and value (as indicated by forward citations) or if local search (narrower recombination) is more likely to lead to cognitive and economic breakthroughs. This approach will allow us to understand if there are any contradictions between the processes leading to novelty and those engendering economic value. We will accomplish this through a test of mediation (Baron & Kenny 1986; Iacobucci *et al* 2007; Zhao *et al* 2010) so that we can examine each of the paths and their joint effects. In testing paths C and C', we replicate prior innovation studies showing the association between recombination and economic value. In testing paths A (and A') and B, we explore the relationship between tension and foundational views in producing novel ideas that should subsequently be associated with economic value.

### **Foundational vs. tension theories and economic value**

Tension assumptions about recombination (Hargadon & Sutton 1997; Weisberg 1999) have dominated management scholarship on innovation. The view is that deep knowledge in a single or small number of domains may lock inventors into one way of thinking and therefore block their ability to generate breakthrough innovations. Local search will only produce incremental innovations and, therefore, to generate breakthroughs, inventors must combine knowledge from distant and diverse sources. In theory, all innovations are based on some sort of recombination. Whenever we use the term recombination, we are referring to the recombination of distant or diverse knowledge, where recombination of local or similar knowledge would be considered local search (e.g., Fleming 2001; Ethiraj & Levinthal 2004).

Drawing on tension view assumptions, researchers investigating the sources of breakthrough patents have identified recombination processes as their source. Their studies have examined a series of recombination measures to show that breakthroughs are the product of combinations of distant and



diverse knowledge. In addition, previous research demonstrates that inventors with experience in recombination and situated in contexts that are conducive to recombination (such as teams and organizations) are more likely to produce breakthroughs. In all cases, these studies examining the sources of breakthroughs use forward citation counts as a measure breakthrough, which captures the economic value of the patent. The relationship they test is represented by path C in Figure 1.

The literature has operationalized recombination in a number of ways. Research has suggested that recombination is most likely to lead to higher citation rates if the knowledge combined is technologically distant. The idea is that combining knowledge from exploratory or long jump search (Gavetti & Levinthal 2000; March 1991) is more likely to produce inventions that break from the existing technological and scientific models and ultimately become highly cited (Phene *et al* 2006, Rosenkopf & Nerkar 2001, Trajtenberg *et al* 1997). Similarly, scholars have argued that highly cited patents are more likely to be combinations of not just distant but also diverse knowledge domains (Hall *et al* 2001), where greater diversity (lower concentration) of knowledge avoids intellectual lock in.

Fleming (2001) has extended these ideas to suggest that if inventors are familiar with the components of an invention and their prior combinations, they will be more able to create new combinations that are valuable. Scholars have also suggested that the degree to which a patent draws on basic scientific knowledge is associated with its future citations. The logic is that science serves as a map for locating innovative combinations (Fleming & Sorenson 2004), and therefore, inventors' embeddedness in science increases the likelihood of finding the most valuable ones (Gittelman & Kogut 2003; Deng *et al* 1999).

Inventors are also seen to be able to recombine ideas better when they collaborate, which, theory suggests, prevents the inertial thinking that any one inventor might experience (Audia & Goncalo 2007). Specifically, top cited patents are more likely to be associated with larger inventive teams due to the greater diversity of viewpoints represented as well as the higher capacity to iterate ideas and select better ones (Singh & Fleming 2010). Teams with greater inventive experience on average will also be more skilled in recombination and therefore better able to create valuable

inventions (Singh & Fleming 2010; Conti *et al* forthcoming). Further, inventors embedded in organizations – especially those that have norms encouraging exploration (Audia & Goncalo 2007) – will be more likely to produce breakthroughs because they are able to recombine a rich amount of knowledge accumulated collectively in the organization (Singh & Fleming 2010). Taking these insights together, we hypothesize that:

*H1a: A patent based on recombination processes is more likely to receive a higher number of citations than a patent produced based on local search (path C in Figure 1).*

To operationalize the general construct of ‘recombination processes,’ we will replicate the variety of measures used by prior scholars mentioned above. Specifically, we will examine the impact of increased distance of knowledge, increased diversity of knowledge, increased familiarity of components and combinations, greater use of science as a map for recombination, invention within teams and organizations that have more resources for recombination, and increased experience of inventors in making recombinations.

While the existing empirical evidence for breakthrough patents has overwhelmingly supported the tension view of recombination, the foundational view would propose an alternative relationship between recombination and measures of economic value. Under this logic, inventors should need to explore a relatively narrow domain in-depth in order to know how to “defy the crowd” and “buy low and sell high” (Sternberg & O’Hara 1999; Sternberg & Lubart 1995). Inventors cannot see new sources of value without understanding what assumptions are behind the existing sources of value, and these insights can only come from focused, local search. This view of creativity is consistent with ecological theories that domain-spanning activities may suffer market penalties due to both deficiencies in production of innovations as well as problems of market reception. Recombinations of distant or diverse knowledge might get in the way of identifying value because they would disperse effort and distract from obtaining the incisive insight that comes from a deep appreciation of one domain (Hannan, Polos & Carroll, 2007). Thus, recombination could prevent the realization of economic value because it produces superficial or incremental work. One might also infer that recombination could

compromise the realization of economic value because market audiences penalize offerings that span categories (Hsu, Kocak & Hannan 2009; Rao, Monin & Durand 2005; Ruef & Patterson 2009; Zuckerman 1999). That is, recombination of inputs could potentially lead to difficulties in classifying the innovation and thus to penalties in the form of fewer citations over time. We therefore offer a competing hypothesis (as represented in Path C') to the recombination model of creativity in generating economic value:

*H1b: A patent based on local search (narrower recombination) is more likely to receive a higher number of citations than a patent produced based on recombination processes (path C' in Figure 1).*

### **Foundational vs. tension theories and novelty**

By calling out novelty as a separate creative output from the generation of economic value, we are able to interrogate existing research that has privileged recombination processes as the source of innovative breakthroughs. Implicit in the arguments made in studies of breakthroughs is the idea that recombination generates novel ideas (path A in Figure 1), which in turn are more likely to be cited as prior art by subsequent patents (path B).

With regard to the connection between novelty and economic value, while the creativity literature makes it clear that not every truly novel idea will become valuable (Amabile 1983; Sternberg & O'Hara 1999; Sternberg 1997), they also indicate that novelty will increase the probability that economic value can be obtained, all else equal. This logic is consistent with the arguments made in the innovation literature on patents discussed above. Of course, it is a requirement of the US Patent Office that every patent be novel to some extent, though some inventions may be 'improvements,' built upon existing technological trajectories, while others may be truly novel, introducing new technological trajectories. Our concern here is with those that meet this latter standard. Thus, we hypothesize:

*H2: Patents that represent cognitive breakthroughs (truly novel ideas) are more likely to receive a higher number of citations than patents that do not (path B in Figure 1).*

As reviewed above, novelty has been portrayed in the innovation literature as an (unmeasured) output of recombination and an input to the creation of economic value (citations). For example,

Trajtenberg *et al* (1997: 29) claim that ‘synthesis of divergent ideas is characteristic of research that is highly original.’ Similarly, Phene *et al* (2006: 370) suggest that, ‘Knowledge that is technologically... distant provides the organization with an opportunity to make novel linkages,’ and Singh and Fleming (2010: 52) claim that, ‘collaboration in the form of team and/or organization affiliation enables more careful and rigorous selection of the best ideas while also increasing the combinatorial opportunities for novelty.’ These arguments are based in the tension view of recombination, which assumes that knowledge and creativity are opposing forces, such that ‘knowledge may provide the basic elements...out of which are constructed new ideas, but in order for these building blocks to be available, the mortar holding the old ideas together must not be too strong’ and too much knowledge of a domain can be habit-forming and inertial (Weisberg 1999, p. 226). Recombination of distant or diverse knowledge can break the habits and inertia.

In introducing a measure of novel ideas, we can make the implicit model in innovation studies explicit: novel ideas – what we can conceptualize as cognitive breakthroughs – are the products of recombination processes:

*H3a: Patents produced through recombination processes are more likely to be cognitive breakthroughs (truly novel ideas) than those that are not (path A in Figure 1).*

The ‘foundational’ view makes the opposite claim (Weisberg 1999). Here, immersion in a particular domain is required in order to produce novelty (Csikszentmihalyi 1996). Local search and narrower recombinations based on deep knowledge in one area enables the identification of anomalies that lead to new insights by exposing the tensions or challenges in the current ways of thinking. This is consistent with the Kuhnian (1962/1996) model in which paradigm shifts are triggered by the accumulation of anomalies. Narrow but deep search leads to truly novel breakthroughs in knowledge because it enables researchers to identify ‘what rules to break’ (Taylor & Greve 2006, p. 726). These findings are also consistent with work at the inventor level of analysis suggesting that specialization is important to push the frontier of knowledge outward as the ‘burden of knowledge’ increases over time (Jones 2009, Agrawal *et al*, 2012, Conti *et al* forthcoming). We therefore offer a competing hypothesis

(as represented in Path A') to the recombination model of creativity:

*H3b: Patents produced through local search (narrower recombination) are more likely to be cognitive breakthroughs (truly novel ideas) than those that are not (path A' in Figure 1).*

If developing truly novel inventions were the only mechanism through which recombination processes would lead to higher economic value, we should expect full mediation of Path C when introducing Paths A and B. However, there are reasons to expect that this may not be the case. Indeed, while the primary mechanism that scholars of innovative breakthroughs theorize is that of novelty, they have not claimed that novelty is the sole driver of economic value (as measured by citations) nor that recombination only serves to generate novel ideas and has no direct effects on the degree of economic value created. However, without a measure of novel ideas, they have not been able to tease apart the effect of novelty from other effects of recombination. For example, combining diverse knowledge domains might enlarge audiences for the innovation or increase the likelihood it will be found by inventors or patent examiners in a search for prior art. Similarly, working in a team might not only enhance novelty through combining the ideas of different members, but could also broaden the network in which innovations would diffuse. To date, scholars have mainly treated these effects as alternative explanations to be controlled for in their analyses so that they can make stronger claims about the implicit relationship between recombination and novelty (e.g., Singh & Fleming 2010).

We do not propose here to test all of the alternative effects of recombination processes and thus would not expect full mediation of Path C when introducing the topic-modeling based measure of novel ideas into the analysis. To demonstrate the distinctive effects of recombination or local search on novelty and economic value (while at the same time taking into account that the generation of novelty is on the path to the eventual realization of economic value), we simply need to find that the combination of Paths A and B is significant. To the extent that we find path C remains significant even when introducing the measure of novel ideas, we would be reinforcing the notion that recombination processes are not only about the generation of novelty or that they support certain parts of the creative process (e.g., selection and retention) and not others (e.g., variation).

## TOPIC MODELING OF PATENT TEXTS: A MEASURE OF NOVEL IDEAS

Crucial to our analysis is the introduction of topic modeling as a way to create a new measure of novelty to contrast with existing citation-based measures of economic value in patents. The intuition behind topic modeling as a method to identify novel ideas is the following: the algorithm uses the co-location of words in a collection of documents to infer the underlying (or latent) topics in those texts and the weight of each topic in each individual document. We can then identify the documents that are the originators of each topic as those early documents with a significant weight in the topic. These originating documents can be seen as cognitive breakthroughs. In our case, because we study patents, we call these topic-originating patents. Because topic modeling is a new method in strategic management, we introduce it first here before delving, in the next section, into the empirical methods and measures of other variables, which are more standard in the field. We explain how our method works and then show how we have implemented it in our sample of fullerene patents in order to construct the measure of novelty.

Our methodological move is to treat the texts of patents as representations of the inventive ideas embodied in them. Bibliometric techniques to understand the evolution of science and technology have a long tradition starting from the pioneering work of de Solla Price (1965a; 1965b). However, most of the work to date has used citation analyses (e.g., Leydesdorff *et al* 1994; Meyer *et al* 2004; Dahlin & Behrens 2005). Text analysis has been much less frequent, and, until recently, the main uses of the texts were counts, factor analyses and co-word analyses of keywords (typically in the titles of papers or patents) (Yoon & Park 2005; Azoulay *et al* 2007; Mogoutov & Kahane 2007, Upham *et al* 2010). With the increasing power of computation and availability of texts in electronic form, scholars are exploring the possibilities of more comprehensive uses of the texts, which would therefore require automated (unsupervised) approaches. In the field of technology, some recent studies have developed text-based techniques to identify overlaps between documents (Gerken & Moehrle 2012, Winston-Smith & Shah 2013).

The study reported here follows these recent trends. It is premised on the idea that studying

language in documents should provide a reading of their cognitive content (Durlauf *et al* 2007, Whorf, 1956). In management studies, this idea has been adapted methodologically to use word counts to represent themes (Huff, 1990; Abrahamson & Hambrick, 1997; Kaplan *et al* 2003). Where the concern is in identifying themes over large numbers of texts, topic modeling – a text analysis technique developed in computer science – offers exciting potential (see Blei 2012 for an overview; and also, Ramage *et al* 2009; McFarland *et al* 2013 for details). The advantage of topic modeling over word counts and keyword analyses is that it allows for polysemy – words can take on different meanings depending on their contexts – and it is inductive – the scholar does not have to specify categories *a priori* but can allow them to emerge from the data.

Thus, we believe topic modeling should be a fruitful approach to measuring interpretations in the emergence of a new technological field (Hall, D. *et al* 2008). For our purposes, we use the texts in the abstracts of patents to understand how different actors describe what the technology is and could be, and then to identify shifts in language that represent the emergence of novel ideas. We describe the specific choices we made regarding the selection of fullerene patent texts below, but first we offer a primer on the topic modeling procedures we use.

### **A primer on topic modeling**

The goal of topic modeling techniques as developed in the computer sciences is unsupervised analysis of text designed both to generate a predictive model to aid search and to provide a representation of the topics in an existing corpus (Hall, D. *et al* 2008; Chang *et al* 2009). We focus on this second goal, and we will use our data to track the emergence of new meanings over time and identify the patents that lead to these shifts in language.

The topic modeling approach we use is based in the Bayesian statistical technique of Latent Dirichlet Allocation (LDA)<sup>1</sup> (Blei *et al* 2003). Topic modeling allows the researcher to uncover automatically themes that are latent in a collection of documents and to identify which composition of

---

<sup>1</sup> The Dirichlet distribution is a ‘distribution over distributions’ that gives the probability of choosing a group of items from a set given that there are multiple states to consider (it is a distribution over multinomials). Blei *et al* (2003) provide more details on LDA and its comparison with other methods such as latent semantic analysis.

themes best accounts for each document. The documents and the words in the documents are observed but the topics, the distribution of topics per document and distribution of words over topics are unobserved and represent a ‘hidden structure’ (Blei 2012). Topic modeling uses the co-occurrence of observed words in different documents to infer this structure. According to Blei (2012, p. 79), ‘This can be thought of as ‘reversing’ the generative process – what is the hidden structure that likely generated the observed collection?’ Computationally, the algorithm identifies the posterior distribution of the unobserved variables in a collection of documents.

This idea is represented schematically in Figure 2. The shaded circle denotes what can be observed ( $\mathbf{w}$ , the words in the documents in the collection). The unshaded circles denote latent (unobservable) variables:  $\mathbf{z}$ , the topic assignment in each document;  $\boldsymbol{\theta}$ , the per-document topic proportions; and  $\boldsymbol{\alpha}$  and  $\boldsymbol{\beta}$ , the parameters of the Dirichlet priors for  $\boldsymbol{\theta}$  and the distribution of topics over words respectively. Each box is a plate where the N plate denotes the words within documents, the D plate denotes the documents within the collection and the T plate denotes the distribution of words over topics. Each word is assumed to be drawn from one of T topics. All topics are used in every document, but exhibit them in different proportions (usually where a few topics are quite important and most are hardly salient). The arrows indicate conditional dependencies between the variables such that assigning topics to words ( $\mathbf{z}$ ) depends on the per-document topic proportions ( $\boldsymbol{\theta}$ ), and the appearance of a word in a document is inferred to be dependent on the distribution of topics over words ( $\boldsymbol{\beta}$ ) and the topics in each document ( $\mathbf{z}$ ).

-- Insert Figure 2 about here --

Given a collection of documents, the topic model algorithm provides two outputs, the first being a list of topics with a vector of words weighted by their importance to the topic, and the second being a list of documents (in our case, patent abstracts) with a vector of topics weighted by their importance to the document. This method allows the researcher to quantify meaning over large numbers of texts and to identify shifts in thought. A feature of this approach – a feature that moves beyond simple word counts or overlaps – is that the same word may have different meanings



depending on its co-occurrence with other words in a document, where the number of topics in a word corresponds with the number of different meanings it has (Chang *et al* 2009). This is particularly important given Kuhn's (1962/1996, p. 205) argument that 'proponents of different theories are like the members of different language-culture communities,' where vocabularies might share many of the same words but the actors attach different meanings to them.

A topic is a multinomial over a set of words, and therefore is not labeled by the algorithm (Blei & Lafferty 2007). Though some scholars have experimented with the automatic labeling of topics, this approach is not reliable enough to have been widely adopted (Mei *et al* 2007). Thus, a further step in using topic models is the labeling of topics based on the words in them, which serves an important function in validating the topics produced by the model as well as generating a label to characterize each topic. As described below, we engaged three nanotechnology experts to label and validate the topics generated by the topic model.

We used the publicly available 'Stanford Topic Modeling Toolbox' developed by the Stanford Natural Language Processing Group and made available in 2009 (Ramage *et al* 2009).<sup>2</sup> The algorithm requires inputs for the two parameters  $\alpha$  (sometimes called 'topic smoothing') and  $\beta$  (sometimes called 'term smoothing'). Values above one lead to more even distributions. Values below one favor more concentrated distributions across fewer topics or words. In order to produce semantically meaningful topics, the Stanford Topic Modeling Toolbox recommends 0.1 for both parameters as a default. As a smaller  $\beta$  results in more fine-grained topics (Griffiths & Steyvers, 2004), we lowered this parameter to 0.01 because we are studying a narrow field of technology. The algorithm thus allocates documents to the fewest topics possible while at the same time assigning a high probability to as few words as possible for each topic.

For computer science applications such as the development of predictive models for text searches, the best-fit model often produces a very large number of topics. However, Chang *et al* (2009)

---

<sup>2</sup> See <http://nlp.stanford.edu/software/tmt/tmt-0.4/> for further details on the toolbox. Another good option is MALLET (<http://mallet.cs.umass.edu/>). Many new implementations are emerging as topic modeling becomes more prevalent, e.g., one can use the R package recently developed by Grun and Hornik (2011).

show that these best-fit models do not produce topics that represent distinct meanings and that smaller numbers of topics make interpretation more feasible. Thus, scholars have found it most useful to constrain the number of topics (where the typical number selected is 100) (Blei & Lafferty 2007; Hall, D. *et al* 2008). Following their lead, we limited the model to 100 topics, which was the maximum number that would still be interpretable by fullerene experts. This provided both statistically and semantically meaningful topics.

### **Sample of fullerene and related patents**

To test the use of topic models to identify shifts in vocabularies, we focused on a single technical domain, that of buckminsterfullerenes (and the chemically related carbon nanotubes). This narrow focus is essential because it allowed us to identify field-level experts to validate the topics generated using the topic-modeling algorithm. Prior studies of the emerging field of nanotechnology have found fullerenes and nanotubes to be a useful site for analysis (Kuusi & Meyer, 2007; Wry, *et al* 2010) because they can be applied in a broad range of potential applications from medicine to electronics to sports and therefore can be conceptualized as general purpose technologies (GPTs) (Bresnahan & Trajtenberg, 1995; Helpman, 1998).<sup>3</sup> They have the chemical formula of C<sub>60</sub> or Carbon 60. Buckminsterfullerenes (also known as fullerenes) were discovered in 1985 by Dr. Richard Smalley, Robert Curl and Harold Kroto (for which they won the Nobel Prize in Chemistry in 1996). Carbon nanotubes are in the fullerene family and their discovery is attributed to Sumio Iijima of NEC Corporation in 1991.

The choice of fullerenes and nanotubes is appropriate for the application of topic modeling because they are subject to substantial patenting over time and are associated with a multiplicity of interpretations. These patents show that inventors envision technologies for revolutionary new applications (e.g., implantable medical devices to control insulin levels for diabetics, more targeted treatments for cancer, structural materials for combat and sports gear, super lightweight batteries and

---

<sup>3</sup> Fullerenes compare quite well in the generality index with other GPT's studied (Hall, B. *et al*, 2001). For example, from 1990 to 2000, the average generality index of fullerene patents is above 0.6 which is nearly 50 percent higher than that reported for patents in computers and communications (the technologies Hall *et al* identify as having the highest generality).

new computing processors that provide quantum leaps in speed and storage capability). Because of this range of potential applications, researchers and managers in universities and firms have broad purview to guide the research and development of the technology in many directions. As a result, their interpretations of what the technology is and how it might be used have consequences for the development and evolution of the technology. Research and development (and ultimately commercialization) resources will be placed in some areas and not others depending on the interpretations and choices these researchers make.

We collected the 2,826 fullerene and nanotube patents granted by the US Patent and Trademark Office (USPTO) through 2008 (stopping in 2008 allows us to collect 5-year forward citation data for all patents in the sample). We identified the population of patents using three separate search techniques. First, we used Derwent's technology classifications to select all patents they identify as pertaining to either of these technologies: B05-U; C05- U; E05-U; E31-U02; L02-H04B; U21-C01T; X12-D02C2D; X12-D07E2A; X12-E03D; X16-E06A1A. Second, the USPTO established a nanotechnology 'cross reference' class (#977) in 2004, which was applied retroactively to all previously-granted patents deemed relevant as well as to new nanotechnology patents. All the patents in subclasses pertaining to fullerenes and nanotubes (977/735-752) were selected. To complement the use of these formal classification systems, we also selected all utility patents with the terms 'fullerene' or 'carbon nanotube' in the title, abstract or claims.

Creating a comprehensive set of patents is imperative because the output of topic modeling depends on the collection of documents used. Figure 2 demonstrates that no individual sampling technique provided a complete picture of patents that could plausibly be associated with fullerenes, and we believe our approach to developing the population of patents in this field compensates for biases created by any one method of classification.

-- Insert Figure 2 about here --

### **Deriving fullerene and nanotube topics**

For each patent, the abstracts from its USPTO document were used. Patents' abstracts are particularly

appropriate for an analysis of shifts in language because they are meant to represent a summary of the novel aspects of the invention. Specifically, the USPTO instructs applicants that, ‘The purpose of the abstract is to enable the [USPTO] and the public generally to determine quickly from a cursory inspection the nature and gist of the technical disclosure,’ where, ‘the form and legal phraseology often used in patent claims... should be avoided’<sup>4</sup> (see also Emma 2006). Further, the USPTO requires that abstracts be 150 words or less, thus assuring that the documents compared in our analysis are of approximately equal size.

In several cases, multiple patents with the same abstract have been granted to protect a single invention. To prevent multiple counting of such texts, we grouped patents with identical abstracts and assignees into patent families. This resulted in 2,384 patent families based on the 2,826 patents (there are 336 families with more than one patent, most of which include only 2 patents, with an average of 2.56 and a maximum of 15). We use the data associated with the chronologically first patent in the family. As is typical practice, we removed ‘stop words’ such as ‘the,’ ‘and,’ ‘that,’ or ‘were’ that do not contribute to the identification of topics. Using the approach described above, we identified 100 separate topics, the probability that each of the words appeared in each topic, and the weight of each topic in each abstract.

To label the topics and validate their usefulness in identifying separate ideas, three nanotechnology experts<sup>5</sup> separately reviewed each of the 100 topics. Based on the list of the top 20 words and their weights as produced by the topic modeling algorithm, we asked each coder to provide a short name to label the topic. We obtained a Krippendorff’s  $\alpha$ , a common measure of inter-rater reliability, of 0.78, which is high given that the coders were not starting from an initial list of codes to apply to each topic. Disagreements for 22 topics were all resolved in joint discussions.<sup>6</sup> A series of

---

<sup>4</sup> First quote from, [http://www.uspto.gov/web/offices/pac/mpep/documents/0600\\_608\\_01\\_b.htm](http://www.uspto.gov/web/offices/pac/mpep/documents/0600_608_01_b.htm) (accessed November 15, 2012). Second quote from the Manual of Patent Examining Procedure (MPEP), Eighth Edition, August 2001, Latest Revision July 2010, Chapter 600 Parts, Form, and Content of Application, 608.01(b) ‘Guidelines for the Preparation of Patent Abstracts’ Section C: ‘Language and Format’

<sup>5</sup> Two post docs and one graduating PhD student, each with experience in research on fullerenes and nanotubes.

<sup>6</sup> Coders found 25 topics to be very general and therefore difficult to label with distinctive codes. This is not surprising as topic models tend to place noisy data into broad or uninterpretable topics (which serves to bolster the coherence of the other topics). We performed all of our statistical analyses omitting these general topics. We found fully consistent results (results available from the

topics focused on production processes such as chemical functionalization of nanotubes, metal catalysts for production, or using a reaction vessel for producing nanotubes. Other topics covered applications into such areas as neural networks, reinforced golf balls, optical devices, batteries, transistors, magnetic memory, recording devices, temperature sensing devices, x-ray devices, DNA detectors or plasma display panels. A third category included topics related to the equipment – primarily scanning probe microscopes – used for visualizing and manipulating nanoscale matter.

Figure 3 shows a sample abstract and the weight of its most important topics. This abstract for patent number 7288970 ‘Integrated nanotube and field effect switching device’ is dominated by topic 24 (*‘Nanotube switching devices and applications’*) with 63 percent weight and topics 54 (*‘Electronic implementations of look up tables’*) and 49 (*‘Field emissions display devices’*) each with 5 percent weight. No other topic is greater than 5 percent weight. The sum of the weights of all the topics for any given patent is 100 percent. Figure 4 provides graphical representations of sample topics with the top 20 words associated with each sized in proportion to their importance.

-- Insert Figures 3 and 4 about here --

The topics as generated from the abstracts of patents do not give us the same information as that captured by patent office classifications. The correlation between categories developed using topic modeling and the USPTO technological classes (using primary topics and primary 3-digit patent class) is 0.22 with a standard deviation of 0.10 (where the average correlation is the average over all the calculated maximum correlation values for each topic with all the patent classes). This is perhaps not surprising. In the case of nanotechnologies, the USPTO did not have a standardized classification system for this field until the introduction of the 977 class in 2004 and, even then, 977 was only a cross-reference class, and therefore would never appear as a primary patent class. However, even when examining the correlation between primary topics and 977 class assignments (for the 305 patents that were assigned a 977 class), we find that it is only slightly higher (0.29). This may be the case because

---

authors). However, since our interest is in introducing a replicable approach using unsupervised analysis of texts, we report the results based on all 100 topics below rather than the results that depend on human intervention by coders.

topics are generated from the writings of inventors (and others who help construct the patent) to describe the nature of the invention, while classifications are assigned by patent examiners using previously established classification systems to facilitate their search for prior art (USPTO, 2005).

### **Identifying topic-originating patents that represent truly novel ideas**

There are many potential analytical uses of the data produced using topic modeling. For the purposes of this study, we focus on the identification of patents that originate novel ideas. To do so, we detect the entry of new topics into the sample. We then select all patents over a threshold weighting for that topic (in our case, 0.2) and appearing in the first 12 months of the topic formation (based on application date). The average number of topic-originating patents using this method is 1.89 per topic for a total of 189. The median is 1. Two topics had more than 10 patents associated with them in the first year; results are robust to their omission. Thus, *topic originating patents* is an indicator variable where a 1 identifies those patents that are over the threshold (though results are robust to the use of a continuous variable where topic originating patents are measured according to the weight of the topic in the patent).

The selection of topic-originating patents is sensitive to the cutoff points we set. For the time frame, we chose 12 months as reasonable estimate of the time for which the knowledge of that invention would not be widespread (where the average lag between application and granting of a patent in our data is 34 months). Thus, any patents applied for in this 12-month window could be considered simultaneous inventions (though, in our regressions we test the use of only the first patents in each topic and get similar results as those reported here). The threshold for topic weight is also an important choice. To identify the appropriate threshold, for each of the 100 topics, we provided our three expert coders with a chronological list of patents (and their abstracts) that had a greater than ten percent weight in the topic. They were asked to identify the first patent chronologically that represented the essence of the topic. The weight that best matched the expert assessments was a 0.20 threshold. Nevertheless, to check the robustness of this threshold, we conducted our statistical analyses with topic-originating patents as identified using thresholds of 0.15 and 0.25. These results are qualitatively

similar to those using the best-fit threshold, with the effect of topic originating patents in our regressions slightly lower (but still significant) for patents using the tighter 0.25 threshold and slightly higher (but also still significant) for patents using the less-strict 0.15 threshold.<sup>7</sup>

Though few other scholars studying breakthrough innovations have attempted to measure novelty directly, Ahuja and Lampert's (2001) study of organizations offers some analogues in their measures of 'pioneering,' 'novel' and 'emerging' technologies. The challenge in comparing our measure of novel ideas to their metrics is that they are all calculated at the firm-level of analysis and therefore are not easily transposed to the invention- (patent-) level of analysis we use in our study. If we were to adjust their measure of 'pioneering' technologies from a count of a firm's patents that have no prior art to a dummy indicating if a single patent has no prior art and their measure of 'emerging' technologies from a count of firm patents that cite prior art that is on average less than 3 years in age to the average age of citations for a single patent, we find that they are weakly but positively correlated with topic originating patents (0.064 and 0.055 respectively). We could not find an equivalent invention-level measure for their variable 'novel technologies' because this is based on the number of new patent classes entered by the firm in the previous 3 years. Nevertheless, this analysis in some small way validates our invention-level measure of novelty.

One might be worried that the identification of topic-originating patents is merely mechanical. Since topic modeling calculates posterior probabilities, the patents we identify as topic-originating could be supposed to exist only because of the many citations to those patents that follow subsequently or because certain topic originating patents are associated with more heavily populated topics. We do not believe this to be the case for several reasons. First, the number of patents per topic differs from topic to topic, ranging from 9 to 44 (mean 23.2, standard deviation 8.7). Second, in examining the 5-

---

<sup>7</sup> We also explored other methodologies for identifying topic-originating patents, such as selecting the first patent to represent a substantial jump in weight relative to prior patents. For example, if we look at patents that represent a 3 standard deviation jump from prior patents and all subsequent patents over that weight in the first year, we obtain a list of 141 patents, of which 73% are the same as those topic-originating patents identified using our .2 threshold. The results using this alternative measure are substantially the same as those we report in the paper. Note, however, that in every case, the selection of a topic-originating patent requires an assumption about a threshold, whether it is a weight of .2 or a number of standard deviations. Therefore, we prefer to use the threshold that has been validated by coders.

year forward citations for each of the topic-originating patents, we see that there is substantial variation around the mean of 22.17, where the minimum is 0 and the maximum is 187. Indeed, only 20 of the 189 topic-originating patents are also in the top 5 percent of cited patents (those that are often considered ‘breakthroughs’). Moreover, the correlation between the topic size and the number of citations to its originating patents is negative and insignificant. When we enter the topic size as a control variable in our regressions, its estimated coefficient is very small, negative and not statistically significant, suggesting that there is no direct positive relationship between the size of the topic and the number of citations that topic’s originating patents receive from follow-on patents. Many of the forward citations are made by non-fullerene patents (and therefore not in our dataset and not used in the text analysis that identified the topics), and for those that are fullerene patents, the citation pattern does not indicate any skewed dependence on topics. Finally, because prior art is assigned based on patent classifications, the low correlation between such classes and the topics would also mitigate any argument of reverse causality. The advantage of topic modeling is precisely that it allows us to identify both heavily and sparsely populated topics. This technique enables the researcher to identify shifts in vocabulary, whether or not many other patents then continue the conversation.

## **MEASURING AND TESTING THE DOUBLE-EDGED SWORD OF RECOMBINATION**

We will use this measure of truly novel ideas (topic-originating patents) as the mediator variable in an analysis of the creative processes producing these novel ideas and subsequent economic value (where citations capture the value of the patent). We describe the dependent and independent variables below and explain how their relationships will be tested using structural equation modeling.

### **Dependent and independent variables**

#### *Dependent variables*

Breakthrough innovations have typically been measured by the number of ‘forward citations’ (prior art citations made to the focal patent by subsequent patents). Higher numbers of citations indicate that a patent represents a breakthrough (Trajtenberg 1990). Because studies of innovative breakthroughs vary as to whether they use a count of forward citations or a dummy variable indicating the ‘breakthroughs’



(for the top tier of cited patents) as the outcome of interest, we examine both dependent variables in our analyses. We examine the 5-year count of forward citations as a dependent variable in negative binomial count models and a dummy variable indicating whether a patent is in the top 5 percent of cited patents in a probit model. The two dependent variables are measured from the grant date in our reported regressions. Our results are robust to the use of a 5-year window since application date.

### *Independent variables*

To test for the competing effects of the tension and foundational views of recombination, we can draw on measures of recombination from prior studies where higher levels of any of the measures would capture recombination (the tension view) and lower levels would capture local search (the foundational view). Here, we replicate a wide series of studies using patents to examine different aspects of breakthrough innovations.

**Technological distance.** To measure the distance of the knowledge recombined, we use the technological distance measure proposed by Trajtenberg *et al* (1997). This measures the distance of the focal patent’s prior art based on USPTO patent classifications as follows:

$$technological\ distance_i = \frac{\sum_j \frac{tech\ distance_{i,j}}{\#\ of\ backward\ citations\ of\ i}}{\#\ of\ backward\ citations\ of\ i}$$

where *tech distance*<sub>*i,j*</sub> is 0 if prior art patents *i* and *j* belong to the same 3-digit technological class, 0.33 if they are in the same 2-digit class, 0.66 if they are in the same 1-digit class, and 1 if they are in different 1-digit classes or have no prior art. The higher the value of this measure, the more distant the knowledge combined.

**Technological diversity.** To capture the breadth of recombination, Hall *et al* (2001) use a Herfindahl index of citation concentration in a measure that represents *technological diversity* (which they termed ‘patent originality’). A high value of diversity (a lower concentration of USPTO patent classes in the prior art cited by a focal patent), the more diverse sources of knowledge it combines. As suggested by Hall (2002), we adjusted the measure to correct for the downward bias associated with patents with few citations to prior art.

$$\widehat{\text{Technological diversity}}_i = \frac{\text{number of patents}_i}{\text{number of patents}_i - 1} (\text{Technological diversity}_i),$$

$$\text{where technological diversity}_i = 1 - \sum_j^{n_i} s_{ij}^2$$

and  $s_{ij}$  denotes the percentage of citations made by patent  $i$  to patents in class  $j$ , out of  $n_i$  patent classes. For patents that cite no prior art, this measure cannot be calculated. Therefore, we set the value to zero and include a dummy (*No prior art*) as a control.

Familiarity of components and combinations. According to Fleming (2001), familiarity with the components of an invention and their prior combinations will enable inventors to create new combinations. Familiarity of components is inferred from how frequently and recently patent subclasses have been used previously by other researchers. This variable – *Ln(component familiarity)* – is measured as the average time-discounted count of all previous usage of the focal patent’s subclasses across all patents listed by the USPTO. Following Fleming’s (2001) formulation:

$$\text{Average component familiarity of patent } i = \frac{\sum_{\text{all subclasses } j \text{ of patent } i} I_{ij}}{\sum_{\text{all subclasses } j \text{ of patent } i} 1}$$

$$\text{where } I_{ij} = \sum_{\substack{\text{all patents } k \text{ granted} \\ \text{before patent } i}} 1\{\text{patent } k \text{ uses subclass } j\} \times e^{-\left(\frac{\text{app.date of patent } i - \text{app.date of patent } k}{\text{time constant of knowledge loss (5 years)}}\right)}$$

Similarly, combination familiarity – *Ln(combination familiarity)* – is measured as the time-discounted count of the previous use of the focal patent’s particular subclass combination across all patents listed by the USPTO:

*cumulative comb. use of patent } i =*

$$\sum_{\substack{\text{all patents } k \text{ granted} \\ \text{before patent } i}} \left[ 1\{\text{patent } k \text{ used same comb. of subclasses as patent } i\} \times e^{-\left(\frac{\text{app.date of patent } i - \text{app.date of patent } k}{\text{time constant of knowledge loss (5 years)}}\right)} \right]$$

On the other hand, Fleming (2001) suggests that too much cumulative use of a combination may mean that it has been exhausted of its potential. We control for this possibility using the variable *Ln(cumulative combination)*, which is the same as combination familiarity but without the time discount.

Science intensity. Another means for inventors to develop a map of the scientific terrain and understand what recombinations are possible is to draw heavily on scientific research. This use of science is often represented as the ‘science intensity’ of a patent. We follow the typical operationalization of this construct as a count of the ‘non-patent references’ listed in focal patent (Gittelman & Kogut 2003; Deng *et al* 1999): *# non-patent references*.

Inventor experience. Inventors with greater experience on average will also be better able to generate recombinations (Singh & Fleming 2010; Conti *et al* forthcoming). This is measured as the average number of previous patents by the inventors of the focal patent, using a log normal transformation to deal with the skewness of the data: *Ln(average experience)*.

Teams and organizations. Collaborations are also seen to produce more recombinations, due to their greater diversity of views and backgrounds. We follow the lead of Singh and Fleming (2010) and measure this as a dummy (*Team*) but in separate analyses test the count of team members and find similar results. Relatedly, if an inventor is embedded in an organization rather than being solo operator, s/he can draw on a greater variety of accumulated knowledge to make recombinations (Audia & Goncalo 2007; Singh & Fleming 2010). This is measured according to Singh and Fleming’s (2010) approach as dummy variable (*Assigned*) indicating the patent was assigned to organization.

Additional controls. We include three other measures as controls because they have been shown to be associated with the forward citations garnered by patents. We control for the total number of patents cited as prior art (*# domestic references*) because it is assumed that patents that cite more will also be cited more (Podolny & Stuart 1995). We also control for the number of claims (*# claims*), because it has been argued that the greater the scope of the patent, the more likely the invention will receive future citations (Singh & Fleming 2010). Finally, we control for *family size*, where the family is the set of patents that contain identical abstracts and assignees and therefore are assumed to represent a cluster of patents around a single invention. We assume that patents in large families will be more likely to receive higher numbers of future citations (this is related to arguments by Cockburn & Henderson 1998; Gittelman & Kogut 2003; Harhoff *et al* 2003, who measure patent families as patents

that are patented in multiple jurisdictions). We include annual time dummies as a simple control for possible time trends.

Table 1 shows the means and standard deviations for the whole sample – the majority of which have similar values to those reported in the studies we cite above – as well as for subsamples of topic-originating patents, highly cited patents vs. all others. Note that, relative to other patents, topic-originating patents have higher numbers of citations but statistically significantly lower values for most of the variables measuring recombination, an initial indication of support for the foundational view of creativity in producing novelty. In contrast, highly cited patents have higher values for most of the recombination measures, in line with prior results reported in the innovation literature based on the tension view. Note also that the measure of novelty is positively correlated with economic value. The correlation table (not reported here for reasons of space) confirms these results.

-- Insert Table 1 about here --

### **Structural equation modeling to test for mediation**

We will use structural equation modeling to examine the double-edge sword of recombination. Where there are more than one independent variable (the case in our analysis), structural equation modeling (SEM) is the recommended approach for testing mediated relationships (Iacobucci *et al* 2007; Zhao *et al* 2010; Cho & Pucik 2005). Positive signs for our recombination variables as tested in paths A, A\*B (the indirect effect of recombination on citations as mediated by novelty) and C in Figure 1 would support the tension view of creativity; negative signs (Paths A', A'\*B, and C') would be evidence for the foundational view. This approach will also allow us to verify prior studies showing a direct, positive association between recombination and citations (Path C). The advantage of SEM relative to running three separate regressions (the traditional approach to testing mediation, according to Baron & Kenny 1986) is that the simultaneous equations control for measurement errors that might lead to under- or over-estimation of mediation effects (Shaver 2005). The indirect effect of one of the independent variables (IV) on the dependent variable (DV) through the mediator can be calculated by multiplying the estimated direct effect of the IV on the mediator (path A) and the estimated direct

effect of the mediator on the DV (path B). As a robustness check, we also conducted a mediation analysis with separate regressions for each path in Figure 1 and found highly consistent results in both the effect size and significance.

The nature of our dependent variables (one is a count and the other binary) and mediator variable (also binary) places additional constraints on the SEM approach. Using a linear model with count and categorical dependent and mediator variables can lead to biased results. We therefore use the Generalized SEM (GSEM) model introduced in Stata 13 that allows generalized linear response functions with count and binary outcomes. We use a negative binomial function for regressions with count outcomes and a probit function for regressions with binary outcomes. Employing a maximum likelihood estimator, GSEM provides consistent, efficient and asymptotically normal estimates for paths A, B and C. We further use nonparametric bootstrapping (with 1,000 replications) to adjust estimates for bias and to estimate the indirect effects ( $A*B$ ), total effects ( $[A*B]+C$ ), their standard errors and their confidence intervals. All the significance levels are determined by the bias-adjusted bootstrap confidence intervals (Mooney & Duval 1993, Efron & Tibshirani 1993).

## **RESULTS**

Tables 2 (for citation counts) and 3 (for breakthroughs in the top 5 percent of citations) report the results of the structural equation models. Column 1 shows the direct effect of the independent variables measuring recombination processes and mediator measuring novel ideas on the dependent variable (paths C and B). Column 2 shows the direct effect of recombination processes on novel ideas (path A). Column 3 identifies mediation effects in the analysis and shows the indirect effect of recombination processes as mediated by novel ideas ( $A*B$ ). Column 4 represents the total effect of the independent variables and mediator on the dependent variables, taking into account the direct and indirect effects ( $[A*B]+C$ ).

-- Insert Tables 2 and 3 about here --

### **Testing path C**

In testing H1a and H1b for path C, we are in essence replicating the prior studies on breakthrough

innovations linking recombination with subsequent citations. This is a first step of a test of mediation but also serves to establish the validity of the dataset used in this current study. Though the samples of the prior studies are vastly different (in terms of numbers of observations, time periods, technological arenas, etc.), we find support for H1a in our fullerene patent dataset in terms of direction and, in most cases, significance of effect for each of the citation-based measures (Model 1 in Tables 2 and 3). That is, as previous studies have found, recombination processes – greater distance and diversity of knowledge, greater use of science as a map for recombination, more familiarity with components and combinations, invention in organizations and teams with greater resources for recombination and greater experience of inventors – are positively associated with the generation of economic value in the form of citations to patents.

For the citation count models in Table 2, we find that more distant (*technological distance*) and more diverse (*technological diversity*) combinations of knowledge have a positive, though sometimes only marginally significant, effect on subsequent citations. Further, the coefficients for inventor familiarity with prior knowledge – *Ln(component familiarity)* and *Ln(combination familiarity)* – are positive and mainly significant. We also find that *Ln(# non-patent references)* – used as a proxy for the science intensity of a patent – is positively and significantly associated with citations. *Ln(average experience)* – used as a proxy for experience in recombination – is also positive. Looking at the organizational factors that might promote recombination, we find that the effects of *team* and *assigned* are unambiguously positive and significant. The various controls also operate mainly as expected. The significance of the effects are attenuated when using the dummy variable for citation-based breakthroughs (in Table 3), but this is likely due to the reduction in variance in the dependent variable and to the much smaller number of observations in our sample relative to other studies that have used this outcome measure.

### **Testing mediation (paths B and A)**

Confirming H2, Model 1 in both Tables 2 and 3 shows that the measure of truly novel ideas (topic-originating patents) is strongly positively associated with subsequent citation rates. This relationship

(for path B) is statistically and economically significant. Looking at Model 1 in Table 2, a topic-originating patent is likely to receive 1.4 times more citations than the average patent. Similarly, looking at Model 1 in Table 3, if a patent is topic-originating, the odds of it becoming an economic breakthrough as measured by citation rates increase by a factor of 1.7. In other words, holding all other variables at their means, the probability of gaining a breakthrough level of citations is 0.072 for topic-originating patents (those representing divergent ideas) compared to 0.024 for other patents. The marginal effect of topic originating patents on the likelihood of becoming a top cited patent is 0.031, which is substantially greater than the marginal effects of any of the other recombination variables.

On the other hand, we do not see the positive effect of recombination on novelty anticipated by tension theories (H3a). First, looking at the direct effect of recombination variables on novel ideas in Model 2, we find that *technological distance* and *technological diversity* are negatively and significantly associated with topic-originating patents. That is, topic-originating patents are not the result of the combination of distant or diverse knowledge. No other recombination variables (except experience) appear to have a significant association with novelty as measured by topic-originating patents. The positive and significant signs for experience on both value and novelty suggest that inventors' previous experience in patenting may increase familiarity with recombination (as suggested by the tension view) or deep knowledge in the field (as suggested by the foundation view), or both. Future research might explore the effect of experience on these two different creative processes.

Further, turning to Model 3, which is the test of mediation from structural equation modeling (Iacobucci *et al* 2007, Zhao *et al* 2010), we do not find the complementary mediation relationship hypothesized in H3a. Instead we find support for H3b, which suggests a competing mediation relationship (for some variables). Breakthrough novel ideas are associated with higher citation rates but recombination processes do not appear to produce that novelty, and in the case of the distance and diversity of knowledge recombinations, pull in the opposite direction (a partial confirmation of the foundational view of creativity as represented in H3b). Because distant and diverse recombinations have a positive direct effect on citations but are negatively associated with cognitive breakthroughs,

their total effect (Model 4) is not statistically significantly different from zero. This is the essence of the double-edged sword of recombination.

These competing mediation effects are statistically significant based on multiple tests. Where the dependent variable is citation counts, estimates associated with paths A and B are calculated using different response functions (probit versus negative binomial). Thus, one might be concerned that the estimates for the total indirect effects ( $A*B$ ) and their standard errors might not be accurate. While the properties of GSEM estimates and the bootstrapping method should take care of this concern, we nevertheless performed a robustness test suggested by Iacobucci (2012) to examine the significance of the indirect effects. Here we find the z-statistic of each indirect effect is significant, consistent with the results from the GSEM method. Furthermore, where the dependent variable is a dummy, we used the method proposed by Kenny (2013, see also, MacKinnon & Dwyer 1993) as a robustness check of our results. In this method, the estimates for paths A and B are scaled to similar levels and then their product is estimated using the delta method. Here, again, we find the indirect effects are significant.

### **Relationship between cognitive novelty and economic value**

We find that patents that are especially novel are also especially valuable. Following the methodology introduced by Rysman and Simcoe (2008) to evaluate patent citation patterns and rates adjusted for confounding factors such as cohort effects, we found that topic-originating patents are more likely to have higher citations than other patents both in their first generation and in their second generation (that is, in patents citing patents that cite the focal patent) (results available from the authors).

On the other hand, this relationship is not perfect. As mentioned above, only 20 of the 189 cognitive breakthroughs (as measured by topic modeling) are also economic breakthroughs (patents in the top 5 percent of 5-year forward citations – there are 109 of these in our dataset). Our method thus highlights the separate but interrelated nature of cognitive and economic breakthroughs. Those patents that represent breakthrough levels of both novelty and value are rare (less than one percent of our sample) but appear to have a greater impact on future innovation than any other kind of invention. As a result, our approach may offer empirical handholds for addressing questions of cumulative research,



especially where, as Scotchmer (1991, p. 39) suggests, a ‘first technology has very little value on its own but is a foundation for second generation technologies’ (see also Furman & Stern 2011).

There are several reasons to believe that generating novel ideas may not automatically lead to achieving breakthrough levels of citations (which represent economic value). A novel idea embodied in a patent still depends on other factors to become known and used, including the reputation and status of the inventors (Merton, 1968; Azoulay, Stuart & Wang, forthcoming), the distribution of the idea in the relevant network (Singh, 2005), match of the invention itself with the current environmental demand (Sorensen & Stuart 2000) and the presence of complementary technologies (Rosenberg 1996). In the absence of such factors, a patent that represents a novel idea may not gain traction. Similarly, not all highly cited patents represent divergent breaks in knowledge. Patents with a broad scope and general claims, patents inside patent thickets (dense networks of patents with overlapping claims), patents that make an original idea more understandable and usable or patents that distribute an idea strategically in a network, all may lead to a high level of citations regardless of whether the patent introduces a truly novel idea or not. By adding a direct measure of novelty, our analysis is a first step in separating out the effects of novelty from these other social dynamics (often associated with recombination processes) that should increase value.

## **DISCUSSION AND CONCLUSION**

Our primary objective for this study was to examine the double-edged sword of recombination in creating innovative breakthroughs. To do this, we look at the effect of recombination on two innovative outputs: novelty and economic value. We introduce a new method – topic modeling – for measuring the novelty of ideas embedded in patent texts by identifying those patents that originate new topics in a body of knowledge. This measure allows us to distinguish inventions that are cognitively novel in the Kuhnian sense – they introduce new language and therefore new ways of thinking – from inventions that are economically valuable (as measured by the subsequent citations they receive). It also enables us to examine contrasting creative processes – those based in either ‘tension’ or ‘foundational’ assumptions – that contribute to novelty and value.

## **Implications for understanding breakthroughs**

Our approach takes seriously the idea put forth by Griliches (1990) and pursued in recent studies (Jaffe *et al* 2000; Alcácer & Gittelman 2006; Alcácer *et al* 2009; Benner & Waldfogel, 2008; Hegde & Sampat, 2009; Tan & Roberts, 2010) that patents should be assessed as historical documents produced by inventors, prosecuted by patent attorneys and evaluated by patent examiners. An implication is that it should be useful to analyze the texts in these patents, which is also consistent with the cognitive turn being made in studies of technology emergence and evolution (Kaplan & Tripsas, 2008). In doing so, we complement existing research on technology evolution, in particular that which draws on patent data to understand the sources of innovation.

The imperfect relationship between topic-originating patents and those that receive high citations may indicate that there are different kinds of ‘breakthroughs,’ those that introduce truly novel knowledge and those that are associated with economic value. Distinguishing between the novel and the valuable (and understanding the sources of each) is quite important for several reasons. Breakthroughs in knowledge mark the potential origins of new technological paradigms. Furthermore, identifying the patents that mark shifts in knowledge may help us understand different mechanisms through which new ideas spread over time and space and explain why some new ideas become the wheels of economic fortune and some simply grind to a halt after a few years (Podolny & Stuart 1995).

By operationalizing the concept of novel ideas implicit in many studies of the sources of innovation, we are able to distinguish processes that produce novelty from those that produce economic value. The contrasting results for the measures of recombination are particularly striking. They suggest that generating new topics require deep immersion in a narrower domain rather than linking to more distant or diverse knowledge. On the other hand, patents that cite prior art from a wide range of patent classes are more applicable in a variety of domains and therefore more likely to be cited in the future. These findings highlight the double-edged sword of recombination based on the tension and foundational models of the role of knowledge in creativity (Weisberg 1999; Taylor & Greve 2006). Theories of recombination are based in the former, while our results on the sources of cognitive

breakthroughs are better explained by the latter: local search to uncover anomalies is more likely to produce breaks in the existing knowledge and language. Identifying breakthroughs in knowledge using topic modeling may help us develop further insights into Kuhn's (1962/1996, p. 62) model of technological change based on, 'the previous awareness of anomaly, the gradual and simultaneous emergence of observational and conceptual recognition, and the consequent change of paradigm categories and procedures.'

As an early foray into the use of a new method, this study is, not surprisingly, constrained by some limitations. Most importantly, topic modeling is sensitive to the corpus of documents selected for the analysis. Because the technique is based in the generation of posterior probabilities, the identification of topics and topic-originating patents will be affected by which documents are included in the analysis. This is in turn affected by which inventions are patented and by which documents the researcher selects to include in the corpus. Patents are an imperfect source of information on new scientific and technological ideas. Not all inventions are patented (Scherer 1983; Griliches 1990). We are therefore surely missing ideas and topics that withered on the vine. This constraint is balanced by the rich bibliometric data that patents provide, which allow the scholar to examine the effects of citations, inventors, assignees, patent classes and the like. We have also addressed potential bias in our sample that pre-established classification systems create by using three different search methodologies to identify patents related to fullerenes.

Further, the use of fullerenes as a context may reduce the generalizability of the findings. The selection of this technological field was useful for an initial test of the topic modeling approach because it is a fairly constrained field and yet, because fullerenes function as a general purpose technology (GPT), offers the potential for many different interpretations of what the technology is or could do. Using a narrow context enabled us to identify subject experts who could validate the topics produced by the computer algorithm (Grimmer & Stewart 2013). Our ability to replicate (in testing path C) the results obtained by other scholars on very different datasets gives us some confidence that our 2,826 fullerene patents behave in similar ways as other technologies, at least along the dimensions

we have tested. To the extent that the effect of recombination in producing cognitive and/or economic breakthroughs in a technological domain depends on idiosyncratic characteristics of that domain, these findings may change accordingly. For example, one can imagine that recombination might be more effective in generating novel breakthroughs in more established technological fields where innovators have access to well-developed technological components to recombine. Future research can shed more light on how underlying characteristics of a technological domain would influence these relationships.

### **Implications for organizations**

Our results are shaped by the possibility that the processes we observe are endogenous to each other. We measure the impact of recombination on two different aspects of innovation: novelty and value. Assuming that individuals and organizations are strategic in setting their goals, they simultaneously decide about how much novelty and value they should pursue in their innovative activities. In other words, the decision to achieve a certain amount of economic value is simultaneously determined with the decision to achieve a certain amount of novelty. As a result, what we observe in our data is the resulting outcome of such simultaneous decisions by individuals and organizations about how much effort to place on recombination. In that sense, while pursuing novelty can lead to high citation rates (as we see in our results), pursuing economic value at the same time can influence the level of effort put by individuals and organizations to achieve novel innovations.

This endogeneity has an empirical implication. The simultaneity and mutual dependency between the two decisions directly influence the number of valuable vs. novel innovations we observe in our sample and also the size and significance of the regression coefficients. One can think of another equilibrium in which organizations would have put much more effort in finding novel innovations, which could change the number of topic originating patents in our sample and consequently the size and significance of the effects we find. In that sense, what we measure in paths A and B is influenced by what we measure in path C and vice versa (and this is precisely why we use the GSEM methodology). While this simultaneity and correlation of outcomes can influence our estimated coefficients, our results nevertheless highlight important contrasting effects from recombination

processes on novel and valuable innovative outcomes.

Thus, this endogeneity also has an organizational implication. The results suggest that an effective innovation strategy needs to bridge between recombination and local search in order to facilitate the transformation of novel ideas into economically valuable ones. By understanding the double-edged sword of recombination in driving novelty and value, organizations can think about how to manage these conflicts. The literature has not yet studied the ways in which tension and foundational views of creativity interact. These views have been positioned as alternatives rather than as two processes operating simultaneously in organizations. Our model might be consistent with a variation-selection-retention view of innovation where variation (novelty) is produced by one set of processes while selection and retention of the most potentially valuable ideas is produced by another. Future research could explore the organizational design implications of the presence of these conflicting effects, potentially across different stages of innovation.

Further, while our study has focused on the invention as the unit and level of analysis, research on inter-organizational knowledge spillovers (Jaffe, 1989; Bernstein & Nadiri, 1989; Audretsch & Feldman, 1996; Chacar & Liberman, 2003) and inter-firm competition (Cockburn & Henderson, 1994; Aghion *et al* 2003) may provide additional explanations for the sources and impacts of breakthroughs in novelty, the study of which could offer a fruitful avenue for future research at the firm-level of analysis.

### **Extensions of topic modeling as a tool in studies of innovation**

In addition to identifying the sources and impacts of breakthroughs, topic modeling may usefully contribute to other areas of research on science and technology. For example, topic modeling can allow us to analyze at a more fine-grained level technological distance, ties and spillovers between firms and other entities. To date, this research has primarily been conducted through cross-citation analyses of the overlaps in USPTO patent classifications amongst the patents of different entities (e.g., Jaffe 1986; Ahuja 2000; Song *et al* 2003) or citations between entities (e.g., Jaffe *et al* 1993; Mowery *et al* 1996; Henderson *et al* 1998). Scholars are increasingly raising concerns about the degree to which patent

classifications are proxies of location in technological space (Benner & Waldfoegel 2008) and about the noisiness of patent citations as measures of knowledge flows (Duguet & MacGarvie 2005; Alcacer & Gittelman 2006; Roach & Cohen, 2013). In the case of nanotechnology, for example, ethnographic research has found that knowledge flows are not fully captured by co-authoring and citations, where exchanging students and experimental materials, commenting on each other's work, or participating in problem-solving workshops were more powerful and frequent mechanisms (Mody 2011). Yet, we have lacked reasonable alternative quantitative measures for knowledge flows (Roach & Cohen 2013).

Topic modeling of patents may provide one solution to augment existing approaches. Because topic models produce a vector of weights of each topic for each patent, there is an opportunity to evaluate the content of ties using topics and the strength of ties using weights. This approach may be a useful complement to patent classes because it tracks the language of the actors rather than the classifications assigned by others. It also adds greater nuance than available in current cross-citation approaches by, first, examining the ideas directly rather than inferring them from citation ties and, second, allowing for the possibility that connections amongst ideas occur even if specific patents are not cited.

Topic modeling, thus, offers a new means of generating inductively classifications of ideas from texts, which may be advantageous as we look beyond patents to other collections of documents. With the burgeoning interest in classification and categorization (e.g., Lounsbury & Rao 2004; Kennedy 2005; Navis & Glynn 2010; Pontikes 2012), topic models can identify themes or frames as they emerge and evolve over time (Ruef & Nag, 2014; DiMaggio *et al*, 2013). This approach has the distinct appeal of dispensing with the requirement to use pre-established categories or to come up with ex-post classification systems. Instead, the data can speak for themselves, thus allowing the researcher to observe paths that fall away as well as paths that become consolidated over time. As such, topic modeling could be vital in understanding the emergence and institutionalization of new fields. We hope that our early foray into the application of topic modeling to social science questions can instigate further explorations in these directions.

## References

- Abrahamson, E. & Hambrick, D. C. 1997. Attentional homogeneity in industries: The effect of discretion. *Journal of Organizational Behavior*, 18: 513-532.
- Aghion, P., Bloom, N., Blundell, R., Griffith, R., & Howitt, P. 2005. Competition and Innovation: An Inverted-U Relationship. *The Quarterly Journal of Economics*, 120(2): 701-728
- Agrawal, A. Goldfarb, A. & Teodoridis, F. 2012. Does Knowledge Accumulation Increase the Returns to Collaboration? Evidence from the Collapse of the Soviet Union. Rotman School Working Paper.
- Ahuja, G. 2000. Collaboration Networks, Structural Holes, and Innovation: A Longitudinal Study. *Administrative Science Quarterly*, 45(3): 425-455.
- Ahuja, G., & Lampert, C.M. 2001. Entrepreneurship in the large corporation: A longitudinal study of how established firms create breakthrough inventions. *Strategic Mgmt J*, 22(6-7): 521-543.
- Ahuja, G., Coff, R. and Lee, P. 2005. Managerial Foresight and Attempted Rent Appropriation: Insider Trading on Knowledge of Imminent Breakthroughs. *Strat Mgmt J*.26(8): 791-808
- Albert, M. B., F. Narin, D. Avery, P. McAllister. 1991. Direct validation of citation counts as indicators of industrially important patents. *Res. Policy* 20 251–259.
- Alcácer, J., Gittelman, M., & Sampat, B. 2009. Applicant and Examiner Citations in U.S. Patents: An Overview and Analysis. *Research Policy*, 38(2): 415-427.
- Alcácer, J. & Gittelman, M. 2006. Patent citations as a measure of knowledge flows: The influence of examiner citations. *Review Of Economics And Statistics*, 88(4): 774-779.
- Amabile, T.M. 1983. The Social-Psychology of Creativity – A Componential Conceptualization. *J Pers Soc Psychol* 45(2) 357-376.
- Amabile, T.M., J. Pillemer. 2012. Perspectives on the Social Psychology of Creativity. *J Creative Behav* 46(1) 3-15.
- Audia, P.G., Goncalo, J.A. 2007. Past Success and Creativity Over Time: A Study of Inventors in the Hard Disk Drive Industry. *Management Sci*, 53(1): 1-15.
- Audretsch, D. B., & Feldman, M. P. 1996. R&D spillovers and the geography of innovation and production. *Am Econ Rev*, 86(3): 630-640
- Azoulay, P., Ding, W & Stuart T. 2007. The Determinants of Faculty Patenting Behavior: Demographics or Opportunities? *J of Economic Behavior & Orgs*, 63(4), pp. 599-623, 2007.
- Azoulay, P., Stuart, T., Wang, Y. forthcoming. Matthew: Effect or Fable? *Management Science*.
- Baron, R.M., Kenny, D.A. 1986. The Moderator-Mediator Variable Distinction in Social Psychological Research: Conceptual, Strategic, and Statistical Considerations. *J Pers Soc Psychol* 51(6) 1173.
- Benner, M., & Waldfogel, J. 2008. Close To You? Bias and Precision in Patent-Based Measures of Technological Proximity. *Research Policy*, 37(9): 1556-1567.
- Bernstein, J. I., & Nadiri, M. I. 1989. Research and development and intra-industry spillovers: An empirical application of dynamic duality. *The Review of Economic Studies*, 56(2): 249-269
- Bessen, J. 2008. The Value of US Patents by Owner and Patent Characteristics. *Research Policy* 37(5) 932-945.
- Blei, D. M. 2012. Probabilistic topic models. *Communications of the ACM*, 55(4), 77-84.
- Blei, D. M. & Lafferty, J. D. 2007. A Correlated Topic Model of Science (Vol 1, Pg 17, 2007). *Annals of Applied Statistics*, 1(2): 634-634.
- Blei, D. M.; Ng, A. Y.; Jordan, M. (2003). Latent Dirichlet allocation. *Journal of Machine Learning Research* 3: pp. 993–1022.
- Bresnahan, T. F. & Trajtenberg, M. 1995. General purpose technologies: Engines of growth? *Journal of Econometrics*, 65(1): 83-108.

- Campbell D. 1960. Blind variation and selective retention in creative thought as in other knowledge processes. *Psychological Review* 67: 380-400.
- Chacar, A. S., & Lieberman, M. B. 2003. Organizing for Technological Innovation in the US Pharmaceutical Industry. *Advances in Strategic Management*, 20: 317-340
- Chang, J., Boyd-Graber, J., Wang, C., Gerrish, S. & Blei, D. 2009. Reading Tea Leaves: How Humans Interpret Topic Models. In *Proceedings of Neural Information Processing Systems 2009*.
- Cho, H.-J., & Pucik V. 2005. Innovativeness, Quality, and Firm Performance. *Strat. Mgmt. J.*, 26: 555-575.
- Cockburn, I., & Henderson, R. 1994. Racing to invest? The dynamics of competition in ethical drug discovery. *Journal of Economics & Management Strategy*, 3(3): 481-519
- Cockburn, I. M., & Henderson, R. M. 1998. Absorptive Capacity, Coauthoring Behavior, and the Organization of Research in Drug Discovery. *The Journal of Industrial Economics*, 46(2): 157-182.
- Conti, R., Gambardella, A., & Mariani, M. Forthcoming. Learning to Be Edison? Individual Inventive Experience and Breakthrough Inventions. *Organization Science*.  
<http://dx.doi.org/10.1287/orsc.2013.0875>
- Csikszentmihalyi, M. 1996. *Creativity: Flow and the Psychology of Discovery and Invention*. HarperCollins Publishers, New York.
- Dahlin, K. B., & Behrens, D. M. 2005. When is an Invention Really Radical? Defining and Measuring Technological Radicalness. *Research Policy*, 34(5): 717-737.
- Deng, Z., Lev, B. & Narin F. 1999. Science and technology as predictor of stock performance. *Financial Analysts Journal*, 53(3) 20–32.
- de Solla Price, D., 1965a. Is technology historically independent of science? A study in statistical historiography. *Technology and Culture* 6, 553–568.
- de Solla Price, D., 1965b. *Little Science, Big Science*. Columbia University Press, New York.
- DiMaggio, P., Nag, M. & Blei, D. 2013. Exploiting Affinities between Topic Modeling and the Sociological Perspective on Culture: Application to Newspaper Coverage of Government Arts Funding in the U.S. *Poetics*, 41(6): 570-606.
- Duguet, E., MacGarvie, M. 2005. How Well do Patent Citations Measure Flows of Technology? Evidence from French Innovation Surveys. *Econ. Innov. New Techn.*, 14(5), 375-393.
- Duriau, V. J., Reger, R. K., & Pfarrer, M. D. 2007. A Content Analysis of the Content Analysis Literature in Organization Studies: Research Themes, Data Sources, and Methodological Refinements. *Organizational Research Methods*, 10(1): 5-34.
- Efron, B., & Tibshirani, R. 1993. *An Introduction to the Bootstrap*. Chapman and Hall: New York.
- Emma, P. 2006. How to Write a Patent. *IEEE Micro* 26(1) 143-144.
- Ethiraj SK, Levinthal D. 2004. Modularity and innovation in complex systems. *Mgmt Sci* 50(2): 159-173.
- Fleming, L. 2001. Recombinant Uncertainty in Technological Search. *Management Sci*, 47(1): 117-132.
- Fleming, L., & Sorenson, O. 2004. Science as a map in technological search. *Strat Mgmt J*, 25(8-9): 909-928.
- Furman, J.L., Stern, S. 2011. Climbing Atop the Shoulders of Giants: The Impact of Institutions on Cumulative Research. *Am Econ Rev* 101(5) 1933-1963.
- Gavetti, G., & Levinthal, D. 2000. Looking Forward and Looking Backward: Cognitive and Experiential Search. *Administrative Science Quarterly*, 45(1): 113-137.
- Gerken, J.M., Moehrle, M.G. 2012. A New Instrument for Technology Monitoring: Novelty in Patents Measured by Semantic Patent Analysis. *Scientometrics*, 91:645-670.



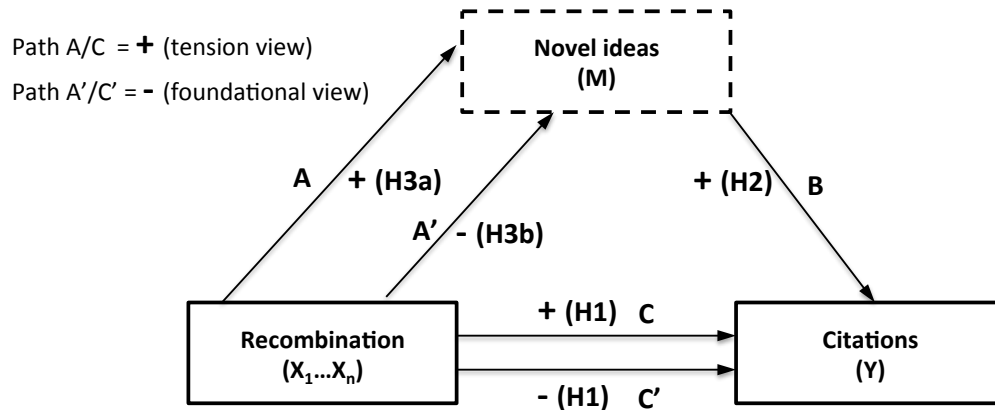
- Gittelman, M. 2008. A note on the value of patents as indicators of innovation: Implications for management research. *The Academy of Management Perspectives*, 22(3), 21-27.
- Gittelman M. & Kogut B. 2003. Does Good Science Lead to Valuable Knowledge? Biotechnology Firms and the Evolutionary Logic of Citation Patterns. *Management Sci*, 49(4): 366-382.
- Griffiths, T. L., & Steyvers, M. 2004. Finding Scientific Topics. *PNAS*, 101: 5228-5235.
- Griliches, Z. 1990. Patent Statistics as Economic Indicators - A Survey. *J of Econ Lit*, 28(4): 1661-1707.
- Grimmer, J. 2010. A Bayesian Hierarchical Topic Model for Political Texts: Measuring Expressed Agendas in Senate Press Releases. *Polit Anal* 18(1) 1-35.
- Grimmer, J. Stewart, B.M. (2013). Text as Data: The Promise and Pitfalls of Automatic Content Analysis Methods for Political Texts. *Political Analysis*. 21(3):267-297.
- Grun, B., Hornik, K. (2011). Topic models: An R Package for Fitting Topic Models. *Journal of Statistical Software*, 40(13): 1-30. URL <http://www.jstatsoft.org/v40/i13/>.
- Guilford JP. 1967. *The nature of human intelligence*. McGraw-Hill: New York.
- Hall, B. H., Jaffe, A., & Trajtenberg, M. 2005. Market Value and Patent Citations. *The RAND Journal of Economics*, 36(1): 16-38.
- Hall, B. H., Jaffe, A., & Trajtenberg, M. 2001. The NBER Patent Citation Data File: Lessons, Insights, and Methodological Tools. NBER Working paper 8498.
- Hall, B. H. 2002. A Note on the Bias of Herfindahl-Type Measures Based on Count Data. In A. Jaffe, M. Trajtenberg, eds. *Patents, Citations, and Innovations*. MIT Press, Cambridge, MA, 454-459
- Hall, D., Jurafsky, D., & Manning, C. D. 2008. Studying the Histories of Ideas Using Topic Models. In *Proceedings of The Conference on Empirical Methods in Natural Language Processing 2008*.
- Hannan MT, Pólos Ls, Carroll G. 2007. *Logics of organization theory: audiences, codes, and ecologies*. Princeton University Press: Princeton, N.J.
- Hargadon A., & Sutton, R. I. 1997. Technology Brokering and Innovation in a Product Development Firm. *Administrative Science Quarterly*, 42(4): 716-749.
- Harhoff, D., Narin, F., Scherer, F. M., & Vopel, K. 1999. Citation Frequency and the Value of Patented Inventions. *The Review of Economics and Statistics*, 81(3): 511-515.
- Harhoff, D., Scherer, F. M., & Vopel, K. 2003. Citations, Family Size, Opposition and the Value of Patent Rights. *Research Policy*, 32(8): 1343-1363
- Hegde, D., & Sampat, B. 2009. Examiner Citations, Applicant Citations, and the Private Value of Patents. *Research Policy*, 105(3): 287-289.
- Helpman, E. 1998. Diffusion of General Purpose Technologies In E. Helpman (Ed.), *General Purpose Technologies and Economic Growth*: 85-117. Cambridge, Mass.: MIT Press.
- Henderson, R., A.B. Jaffe, M. Trajtenberg. 1998. Universities as a Source of Commercial Technology: A Detailed Analysis of University Patenting, 1965-1988. *Rev of Econ & Stats* 80(1) 119-127.
- Hsu G, Kocak O, Hannan MT. 2009. Multiple Category Memberships in Markets: An Integrative Theory and Two Empirical Tests. *American Sociological Review* 74(1): 150-169.
- Huff, A. S. 1990. *Mapping strategic thought*. Chichester, New York: John Wiley and Sons.
- Iacobucci, D. 2012. Mediation analysis and categorical variables: the final frontier. *Journal of Consumer Psychology*, 22(4): 582-594.
- Iacobucci, D., Saldanha, N., Deng, X. 2007. A Meditation on Mediation: Evidence that Structural Equations Models Perform Better than Regressions. *Journal of Consumer Psychology*. 17(2), 140-154.
- Jaffe, A. B. 1986. Technological Opportunity and Spillovers of Research-and-Development - Evidence from Firms Patents, Profits, and Market Value. *Am Econ Rev*, 76(5): 984-1001.
- Jaffe, A. B. 1989. Real effects of academic research. *Am Econ Rev*, 79(5): 957-970

- Jaffe, A. B., Trajtenberg, M., & Fogarty, M. 2000. Knowledge Spillovers and Patent Citations: Evidence from a Survey of Inventors. *Am Econ Rev Papers and Proceedings*, 90(2): 215-218.
- Jaffe, A.B., M. Trajtenberg, R. Henderson. 1993. Geographic Localization of Knowledge Spillovers as Evidenced by Patent Citations. *The Quarterly Journal of Economics* 108(3) 577-598.
- Jones, Benjamin F. 2009. The Burden of Knowledge and the Death of the Renaissance Man: Is Innovation Getting Harder? *Review of Economic Studies*. 76(1), 283-317.
- Kaplan, S., Murray, F., & Henderson, R. 2003. Discontinuities and Senior Management: Assessing the Role of Recognition in Pharmaceutical Firm Response to Biotechnology. *Industrial and Corporate Change*, 12(2): 203-233.
- Kaplan, S. & Tripsas, M. 2008. Thinking about Technology: Applying a Cognitive Lens to Technical Change. *Research Policy*, 37(5), 790-805.
- Kennedy, M.T. 2005. Behind the One-Way Mirror: Refraction in the Construction of Product Market Categories *Poetics* 33(3-4) 201-226.
- Kenny, D. 2013. Mediation with dichotomous outcomes. Working paper, April 19, 2013, available at: <http://davidakenny.net/doc/dichmed.pdf>.
- Kuhn, T.S. 1962/1996. *The Structure of Scientific Revolutions*. Chicago: University of Chicago Press, (3rd ed.)
- Kuusi, O., Meyer, M. 2007. Anticipating technological breakthroughs: Using bibliometric coupling to explore the nanotubes paradigm. *Scientometrics*, 70(3), 759-777.
- Lanjouw, J.O., Schankerman, M. 2004. Patent Quality and Research Productivity: Measuring Innovation with Multiple Indicators. *Econ J* 114(495) 441-465.
- Leydesdorff, L. Cozzens, S., Van den Besselaar, P. 1994. Tracking areas of strategic importance using scientometric journal mappings, *Research Policy*, 23 (2), pg. 217-229.
- Lounsbury, M. & Rao, H. 2004. Sources of durability and change in market classifications: A study of the reconstitution of product categories in the American mutual fund industry, 1944-1985. *Social Forces*, 82(3): 969-999.
- MacKinnon, D.P. & Dwyer, J.H. 1993. Estimating mediated effects in prevention studies. *Evaluation Review*, 17: 144-158.
- March, J. G. 1991. Exploration and Exploitation in Organizational Learning. *Organ Sci*, 2(1): 71-87.
- McFarland, D.A., Ramage, D., Chuang, J., Heer, J., Manning, C.D., & Jurafsky, D. 2013. Differentiating language usage through topic models. *Poetics*. 41(6): 607-625.
- Mei, Q., Shen, X. & Zhai, C. 2007. Automatic Labeling of Multinomial Topic Models. In *Proceedings of The 13th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*.
- Merton, R.K. 1968. The Matthew Effect in Science. *Science* 159(3810): 56-63.
- Meyer, M., Pereira, T.S., Persson, O., Granstrand, O. 2004. The scientometric world of Keith Pavitt, *Research Policy*, 33 (9), pg. 1405-1417.
- Mody, C.C. 2011. *Instrumental Community: Probe Microscopy and the Path to Nanotechnology*. MIT Press, Cambridge, MA.
- Mogoutov, A. Kahane, B. 2007. Data search strategy for science and technology emergence: A scalable and evolutionary query for nanotechnology tracking, *Research Policy*, 36 (6), pg. 893-903
- Mooney, C.Z., & Duval, R.D., 1993. Bootstrapping: A nonparametric approach to statistical inference. No. 94-95. Newbury Park, CA: Sage.
- Mowery, D. C., Oxley, J. E., & Silverman, B. S. 1996. Strategic Alliances and Interfirm Knowledge Transfer. *Strategic Management Journal*, 17: 77-91.
- Natale, F., Fiore, G., & Hofherr, J. 2012. Mapping the research on aquaculture. A bibliometric analysis of aquaculture literature. *Scientometrics*, 90(3): 983-999.

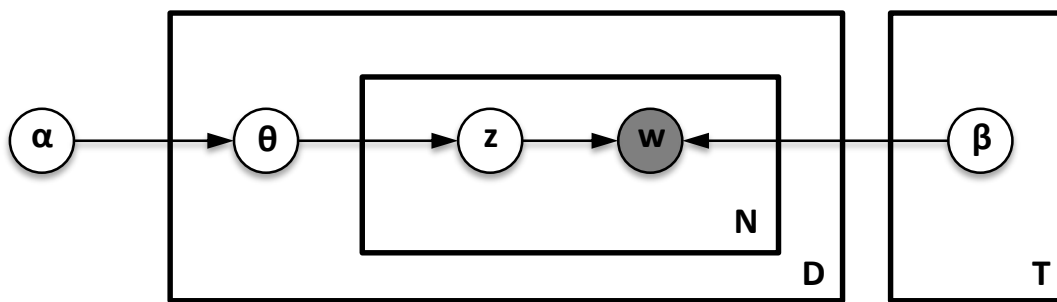
- Navis, C. & Glynn, M. A. 2010. How New Market Categories Emerge: Temporal Dynamics of Legitimacy, Identity, and Entrepreneurship in Satellite Radio, 1990-2005. *Administrative Science Quarterly*, 55(3): 439-471.
- Phene, A., Fladmoe-Lindquist, K., Marsh, L. 2006. Breakthrough Innovations in the U.S. Biotechnology Industry: The Effects of Technological Space and Geographic Origin. *Strategic Management Journal*, 27: 369-388.
- Podolny, J. M., & Stuart, T. E. 1995. A Role-Based Ecology of Technological Change. *The American Journal of Sociology*, 100(5): 1224-1260.
- Pontikes, E.G. 2012. Two Sides of the Same Coin: How Ambiguous Classification Affects Multiple Audiences' Evaluations. *Administrative Science Quarterly* 57(1) 81-118.
- Quinn, K.M., Monroe, B.L., Colaresi, M., Crespin, M.H., Radev, D.R. 2010. How to Analyze Political Attention with Minimal Assumptions and Costs. *Am J Polit Sci* 54(1) 209-228.
- Ramage, D., Rosen, E., Chuang, J., Manning, C., & McFarland, D. A. 2009. Topic Modeling for the Social Sciences. NIPS 2009 Workshop on Applications for Topic Models: Text and Beyond
- Rao H, Monin P, Durand R. 2005. Border crossing: Bricolage and the erosion of categorical boundaries in French gastronomy. *American Sociological Review* 70(6): 968-991.
- Roach, M., Cohen, W.M. 2013. Lens or Prism? Patent Citations as a Measure of Knowledge Flows from Public Research. *Management Sci.* 59(2), 504-525.
- Rosenberg, Nathan. 1996. "Uncertainty and Technological Change." In Landau, Ralph, Timothy Taylor, & Gavin Wright, eds., *The Mosaic of Economic Growth*, Stanford: Stanford University Press, pp. 334-353.
- Rosenkopf, L., Nerkar, A. 2001. Beyond Local Search: Boundary-Spanning, Exploration and Impact in the Optical Disc Industry. *Strategic Management Journal*, 22: 287-306.
- Rothaermel, F.T. & Thursby, M. 2007. The nanotech versus biotech revolution: Sources of Productivity in Incumbent Firm Research. *Research Policy*, (36): 832-849.
- Ruef, M., Nag, M. 2014. "The Classification of Organizational Forms: Theory and Application to the Field of Higher Education," in M.L. Stevens and M.W. Kirst (eds.), *Remaking College: The Changing Ecology of Higher Education*. Stanford, CA: Stanford University Press.
- Ruef M, Patterson K. 2009. Credit and Classification: The Impact of Industry Boundaries in Nineteenth-century America. *Administrative Science Quarterly* 54(3): 486-520.
- Rysman, M., & Simcoe, T. 2008. Patents and the Performance of Voluntary Standard-Setting Organizations. *Management Sci*, 54(11): 1920-1934
- Scotchmer, S. 1991. Standing on the Shoulders of Giants - Cumulative Research and the Patent-Law. *J Econ Perspect* 5(1) 29-41.
- Scherer, F.M. 1983. The propensity to patent. *International J. of Industrial Organization* 1(1), 107-128.
- Schumpeter, J. A. 1939. *Business cycles*. New York: McGraw-Hill.
- Shaver, J.M. 2005. Testing for Mediating Variables in Management Research: Concerns, Implications, and Alternative Strategies. *Journal of Management* 31(3):330-353
- Simonton DK. 1999. Creativity as blind variation and selective retention: Is the creative process Darwinian? *Psychological Inquiry* 10(4): 309-328.
- Singh, J. 2005. Collaborative Networks as Determinants of Knowledge Diffusion Patterns. *Management Science*, 51(5): 756-770
- Singh J. & Fleming, L. 2010. Lone inventors as sources of breakthroughs: Myth or reality? *Management Sci*, 56(1): 41-56.
- Song, J., Almeida, P., & Wu, G. 2003. Learning-by-Hiring: When Is Mobility More Likely to Facilitate Interfirm Knowledge Transfer? *Management Sci*, 49(4): 351-365.

- Sørensen, J.B., Stuart, T.E. 2000. Aging, Obsolescence, and Organizational Innovation, *Administrative Science Quarterly*, 45(1): 81-112
- Sternberg R.J. 1997. *Thinking styles*. Cambridge University Press: Cambridge, U.K.
- Sternberg R.J., Lubart T.I. 1995. *Defying the crowd: cultivating creativity in a culture of conformity*. Free Press: New York, N.Y.
- Sternberg, R.J., L.A. O'Hara. 1999. Creativity and Intelligence. R.J. Sternberg, ed. *Handbook of Creativity*. Cambridge University Press, Cambridge, U.K.; New York, 251-272.
- Steyvers, M., Griffiths, T. Probabilistic topic models. *Latent Semantic Analysis: A Road to Meaning*. T. Landauer, D. McNamara, S. Dennis, and W. Kintsch, eds. Lawrence Erlbaum, 2006.
- Tan, D., Roberts, P.W. 2010. Categorical coherence, classification volatility and examiner-added citations. *Research Policy*. 39(1) 89-102.
- Taylor, A., Greve, H.R. 2006. Superman or the Fantastic Four? Knowledge Combination and Experience in Innovative Teams. *Acad Manage J* 49(4) 723-740.
- Trajtenberg, M. 1990. A Penny for Quotes: Patent Citations and the Value of Innovations. *The RAND Journal of Economics*, 21(1): 172-187
- Trajtenberg, M., Henderson, R., & Jaffe, A. 1997. University versus Corporate Patents: A Window on the Basicness of Invention. *Economics of Innovation and New Technology*, 5(1): 19-50.
- Upham, S.P., Rosenkopf, L., Ungar, L.H. 2010. Innovating Knowledge Communities an Analysis of Group Collaboration and Competition in Science and Technology. *Scientometrics* 83(2) 525-554.
- US Patent and Trademark Office 2005. *Handbook of Classification*. Washington, DC: United States Government Printing Office.
- Weisberg, R.W. 1999. Creativity and Knowledge: A Challenge to Theories. R.J. Sternberg, ed. *Handbook of Creativity*. Cambridge University Press, Cambridge, U.K.; New York, 226-250.
- Whorf, B. L. 1956. Science and linguistics. In J. B. Carroll (Ed.), *Language, thought, and reality: Selected writings of Benjamin Lee Whorf*: 207-219. Cambridge, MA: MIT Press.
- Winston-Smith, S., Shah, S.K. 2013. Do Innovative Users Generate More Useful Insights? An Analysis of Corporate Venture Capital Investments in the Medical Device Industry. *Strategic Entrepreneurship Journal*. 7(2): 151-167.
- Wry, T., Greenwood, R., Jennings, P. D., & Lounsbury, M. 2010. Institutional Sources of Technological Knowledge: A Community Perspective on Nanotechnology Emergence. *Research in the Sociology of Organizations*, 29: 149-176.
- Yoon, B. Park, Y. 2005. A systematic approach for identifying technology opportunities: Keyword-based morphology analysis. *Technological Forecasting and Social Change*, 72(2), 145-160.
- Zhao, X., Lynch, Jr., J.G., & Chen, Q. 2010. Reconsidering Baron and Kenny: Myths and Truths about Mediation Analysis. *J Consumer Res*, 37: 197-206.
- Zuckerman EW. 1999. The categorical imperative: Securities analysts and the illegitimacy discount. *American Journal of Sociology* 104(5): 1398-1438.

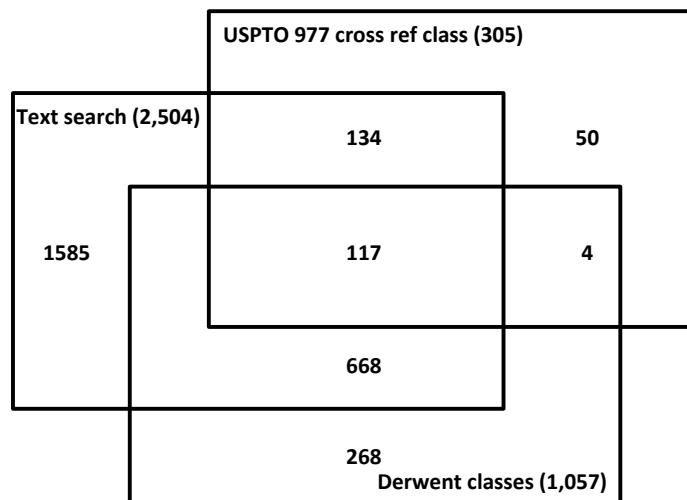
**Figure 1: The double-edged sword of recombination**



**Figure 2: Latent Dirichlet allocation in topic modeling**



**Figure 3: Sample of fullerene and nanotube patents**



Total population = 2,826 patents

## Figure 4: Example coded abstract

**Patent Number:** US7288970

**Title:** Integrated nanotube and field effect switching device

**Inventors:** Bertin, Claude L.

**Assignee:** Nantero, Inc.

**Application Date:** January 2005, **Issue Date:** October 2007

### Abstract

Hybrid<sup>57</sup> switching<sup>24</sup> devices<sup>32</sup> integrate<sup>54</sup> nanotube<sup>56</sup> switching<sup>24</sup> elements<sup>25</sup> with field<sup>49</sup> effect<sup>52</sup> devices<sup>32</sup>, such as NFETs and PFETs. A switching<sup>24</sup> device<sup>93</sup> forms<sup>46</sup> and unforms<sup>24</sup> a conductive<sup>59</sup> channel<sup>52</sup> from the signal<sup>24</sup> input<sup>24</sup> to the output<sup>24</sup> subject<sup>22</sup> to the relative<sup>86</sup> state of the control<sup>24</sup> input<sup>24</sup>. In embodiments<sup>94</sup> of the invention, the conductive<sup>59</sup> channel<sup>52</sup> includes a nanotube<sup>56</sup> channel<sup>52</sup> element<sup>25</sup> and a field<sup>49</sup> modulatable<sup>58</sup> semiconductor<sup>46</sup> channel<sup>52</sup> element<sup>25</sup>. The switching<sup>24</sup> device<sup>93</sup> may include a nanotube<sup>56</sup> switching<sup>24</sup> element<sup>25</sup> and a field<sup>49</sup> effect<sup>52</sup> device<sup>93</sup> electrically<sup>59</sup> disposed<sup>42</sup> in series<sup>24</sup>. According to one aspect<sup>33</sup> of the invention, an integrated<sup>54</sup> switching<sup>24</sup> device<sup>93</sup> is a four-terminal<sup>29</sup> device<sup>93</sup> with a signal<sup>24</sup> input<sup>24</sup> terminal<sup>42</sup>, a control<sup>24</sup> input<sup>24</sup> terminal<sup>29</sup>, a second input<sup>24</sup> terminal<sup>29</sup>, and an output<sup>24</sup> terminal<sup>29</sup>. The devices<sup>32</sup> may be non-volatile<sup>24</sup>. The devices<sup>32</sup> can form<sup>24</sup> the basis for a hybrid<sup>57</sup> NT-FET<sup>52</sup> logic<sup>54</sup> family and can be used to implement<sup>82</sup> any Boolean logic<sup>54</sup> circuit<sup>54</sup>.

Topic 24 (Nanotube switching devices and applications): 63%

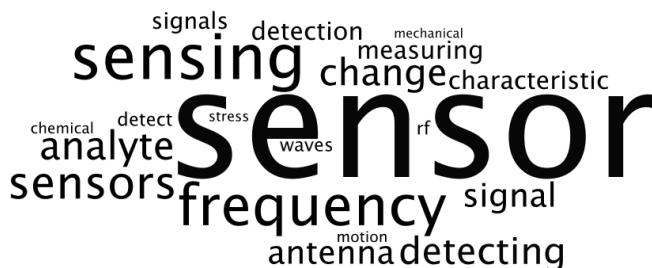
Topic 54 (Electronic implementations of look-up-tables): 5%

Topic 49 (Field emissions display devices): 5%

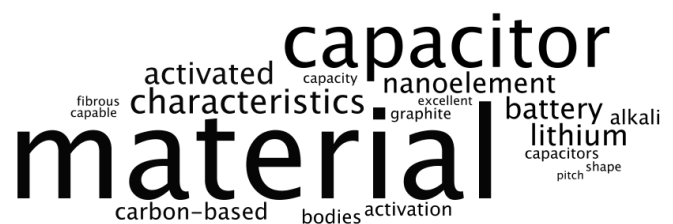
The rest: less than 5% each for a total of 27%

## Figure 5: Graphical representation of words in example topics (date of topic-originating patent)\*

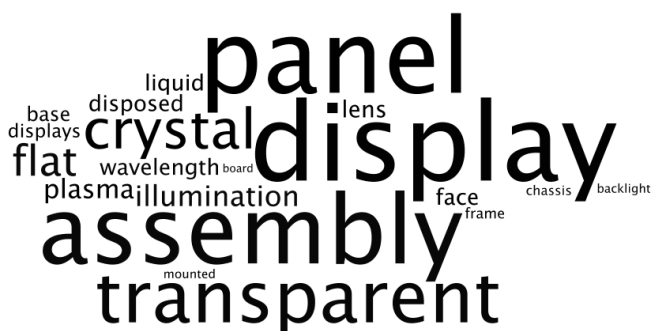
Topic 14: Sensors and detectors  
(1993)



Topic 60: Application to batteries and charge storage devices  
(1994)



Topic 58: Application to plasma display panels  
(1995)



Topic 24: Nanotube switching devices  
(1999)



\* Size of words based on importance in topic

**Table 1: Descriptive statistics, total sample, topic originating patents, top cited patents, and all others**

Means and standard deviations, t-test

	All	Topic originating patent=1	Topic originating patent=0	Difference	Top 5% cited patent=1	Top 5% cited patent=0	Difference
Breakthroughs (top 5% cited)	0.048 (0.213)	0.111 (0.315)	0.042 (0.202)	0.069** (p=0.000)	1.000 (0.000)	0.000 (0.000)	1.000** (p=0.000)
# forward citations (5-yr)	13.507 (21.678)	22.173 (30.958)	12.767 (20.536)	9.406** (p=0.000)	91.982 (29.963)	9.559 (11.106)	82.422** (p=0.000)
Topic-originating patent	0.079 (0.270)	1.000 (0.000)	0.000 (0.000)	1.000** (p=0.000)	0.183 (0.389)	0.074 (0.261)	0.110** (p=0.000)
Technological distance	0.405 (0.328)	0.366 (0.331)	0.409 (0.327)	-0.043+ (p=0.094)	0.474 (0.325)	0.402 (0.328)	0.072* (p=0.026)
Technological diversity	0.736 (0.323)	0.681 (0.353)	0.741 (0.320)	-0.059* (p=0.018)	0.758 (0.317)	0.735 (0.323)	0.023 (p=0.466)
Ln(component familiarity)	4.661 (1.001)	4.337 (1.090)	4.689 (0.989)	-0.353** (p=0.000)	4.735 (0.834)	4.658 (1.009)	0.078 (p=0.430)
Ln(combination familiarity)	0.306 (0.773)	0.264 (0.709)	0.309 (0.778)	-0.045 (p=0.448)	0.225 (0.687)	0.310 (0.776)	-0.084 (p=0.265)
Ln(# non-patent references)	1.498 (1.301)	1.641 (1.341)	1.486 (1.297)	0.156 (p=0.124)	2.360 (1.419)	1.455 (1.280)	0.905** (p=0.000)
Ln(average experience)	2.101 (1.171)	1.961 (1.272)	2.113 (1.162)	-0.152+ (p=0.096)	2.158 (1.197)	2.098 (1.170)	0.059 (p=0.605)
Team	0.791 (0.406)	0.722 (0.449)	0.797 (0.402)	-0.075* (p=0.017)	0.899 (0.303)	0.786 (0.410)	0.113** (p=0.005)
Assigned	0.889 (0.314)	0.889 (0.315)	0.889 (0.314)	0.000 (p=1.000)	1.000 (0.000)	0.883 (0.321)	0.117** (p=0.000)
No prior art	0.048 (0.214)	0.094 (0.293)	0.044 (0.205)	0.051** (p=0.002)	0.027 (0.164)	0.049 (0.216)	-0.021 (p=0.308)
Ln(cumulative combination)	0.409 (0.965)	0.370 (0.871)	0.412 (0.973)	-0.041 (p=0.581)	0.287 (0.819)	0.415 (0.971)	-0.127 (p=0.179)
Ln(# prior art patents)	2.123 (1.066)	1.896 (1.198)	2.142 (1.052)	-0.246** (p=0.003)	2.413 (1.238)	2.108 (1.055)	0.304** (p=0.004)
Ln(# claims)	2.835 (0.764)	2.725 (0.816)	2.844 (0.759)	-0.119* (p=0.045)	3.133 (0.749)	2.820 (0.762)	0.313** (p=0.000)
Ln(Family size)	0.126 (0.331)	0.201 (0.390)	0.120 (0.325)	0.081** (p=0.002)	0.386 (0.547)	0.113 (0.311)	0.273** (p=0.000)

**Table 2: Tests of mediation (dv=citation counts, 5-year window since grant date) (1991-2005)**  
 Generalized structural equation model using bootstrapping (1000 repetitions), bias-corrected coefficients and robust standard errors

	Direct effect on citation counts (Paths C and B) (1)	Direct effect on topic-originating patents (Path A) (2)	Indirect on citation counts (Paths A *B) (3)	Total effect on citation counts ([A * B] + C) (4)
Topic-originating patent	0.330** (0.094)			0.330** (0.094)
<i>Measures of recombination:</i>				
Technological distance	0.211* (0.093)	-0.391* (0.184)	-0.129* (0.073)	0.082 (0.119)
Technological diversity	0.065 (0.101)	-0.486** (0.167)	-0.158** (0.067)	-0.093 (0.118)
Ln(component familiarity)	0.077** (0.029)	0.028 (0.054)	0.009 (0.018)	0.086** (0.034)
Ln(combination familiarity)	0.286 (0.213)	-0.198 (0.328)	-0.068 (0.115)	0.217 (0.241)
Ln(# non-patent references )	0.163** (0.023)	-0.024 (0.041)	-0.007 (0.014)	0.156** (0.026)
Ln(average experience)	0.070* (0.028)	0.122** (0.042)	0.040** (0.018)	0.110** (0.033)
Team	0.226** (0.073)	-0.134 (0.116)	-0.044 (0.041)	0.183* (0.084)
Assigned	0.385** (0.089)	0.020 (0.155)	0.008 (0.054)	0.393** (0.103)
<i>Controls:</i>				
No prior art	0.255 (0.202)	0.491+ (0.279)	0.159+ (0.103)	0.413+ (0.224)
Ln(cumulative combination)	-0.275+ (0.162)	0.192 (0.252)	0.065 (0.090)	-0.209 (0.186)
Ln(# prior art patents)	0.090** (0.032)	0.249** (0.067)	0.082** (0.031)	0.172** (0.044)
Ln(# claims)	0.211** (0.036)	-0.035 (0.064)	-0.012 (0.022)	0.199** (0.042)
Ln(Family size)	0.414** (0.088)	0.113 (0.121)	0.037 (0.042)	0.451** (0.096)
Constant	0.106	0.272	0.089	0.195
Year fixed effects	Yes	Yes	Yes	Yes
Observations	2276	2276	2276	2276

\*\* p<0.01, \* p<0.05, + p<0.10



**Table 3: Tests of mediation (dv=dummy for citation-based breakthroughs, top 5%, 5-year window since grant date) (1991-2005)**

Generalized structural equation model using bootstrapping (1000 repetitions), bias-corrected coefficients and robust standard errors

	Direct effect on on citation-based breakthroughs (Paths C and B) (1)	Direct effect on topic-originating patents (Path A) (2)	Indirect on citation-based breakthroughs (Paths A *B) (3)	Total effect on citation-based breakthroughs ([A * B] + C) (4)
Topic-originating patent	0.538* (0.188)			0.538* (0.188)
<i>Measures of recombination:</i>				
Technological distance	0.044 (0.189)	-0.391* (0.184)	-0.209* (0.126)	-0.165 (0.227)
Technological diversity	0.118 (0.211)	-0.486** (0.167)	-0.257* (0.125)	-0.138 (0.232)
Ln(component familiarity)	0.103 (0.058)	0.028 (0.054)	0.015 (0.031)	0.118 (0.065)
Ln(combination familiarity)	0.564 (0.508)	-0.198 (0.328)	-0.111 (0.196)	0.452 (0.535)
Ln(# non-patent references)	0.247** (0.047)	-0.024 (0.041)	-0.011 (0.023)	0.235** (0.052)
Ln(average experience)	0.079 (0.054)	0.122** (0.042)	0.065* (0.032)	0.145* (0.062)
Team	0.320* (0.175)	-0.134 (0.116)	-0.072 (0.070)	0.248 (0.190)
Assigned	4.522** (0.265)	0.020 (0.155)	0.011 (0.091)	4.533** (0.278)
<i>Controls:</i>				
No prior art	-0.462 (1.093)	0.491+ (0.279)	0.261+ (0.180)	-0.202 (1.119)
Ln(cumulative combination)	-0.519 (0.413)	0.192 (0.252)	0.105 (0.151)	-0.413 (0.444)
Ln(# prior art patents)	-0.062 (0.066)	0.249** (0.067)	0.132* (0.057)	0.070 (0.077)
Ln(# claims)	0.200** (0.072)	-0.035 (0.064)	-0.019 (0.037)	0.180* (0.078)
Ln(Family size)	0.550** (0.128)	0.113 (0.121)	0.061 (0.072)	0.611** (0.144)
Constant	-12.737** (0.650)	0.272 (0.404)	0.142 (0.229)	-12.594** (0.667)
Year fixed effects	Yes	Yes	Yes	Yes
Observations	2276	2276	2276	2276

\*\* p<0.01, \* p<0.05, + p<0.10