

# Measuring Technological Innovation over the Long Run\*

Bryan Kelly<sup>†</sup>   Dimitris Papanikolaou<sup>‡</sup>   Amit Seru<sup>§</sup>   Matt Taddy<sup>¶</sup>

January 2020

## Abstract

We use textual analysis of high-dimensional data from patent documents to create new indicators of technological innovation. We identify important patents based on textual similarity of a given patent to previous and subsequent work: these patents are distinct from previous work but are related to subsequent innovations. Our importance indicators correlate with existing measures of patent quality but also provide complementary information. We identify breakthrough innovations as the most important patents—those in the right tail of our measure—and construct time-series indices of technological change at the aggregate and sectoral level. Our technology indices capture the evolution of technological waves over a long time span (1840 to the present) and cover innovation by private and public firms, as well as non-profit organizations and the US government. Advances in electricity and transportation drive the index in the 1880s; chemicals and electricity in the 1920s and 1930s; and computers and communication in the post-1980s.

---

\*We thank Pierre Azoulay, Nicholas Bloom, Diego Comin, Carola Frydman, Erik Hurst, Pete Klenow, Kyle Jensen, Chad Jones, Chad Syverson, and seminar participants at AQR, Harvard, and the NBER Summer Institute for valuable comments and discussions. We are grateful to Kinbert Chou, Inyoung Choi, Jinpu Yang and Jiaheng Yu for excellent research assistance and to Enrico Berkes and Cagri Akkoyun for sharing their data.

<sup>†</sup>Yale School of Management and NBER

<sup>‡</sup>Kellogg School of Management and NBER

<sup>§</sup>Stanford GSB, Hoover Institution, and NBER

<sup>¶</sup>Amazon

Over the last two centuries, real output per capita in the United States has increased substantially more than the growth of inputs to production, such as the number of hours worked or the amount of capital used. Thus, much of economic growth is attributed to improvements in productivity, which however varies significantly over time and across sectors. Models of endogenous growth ascribe most of these movements to fluctuations in the rate of technological progress. However, both this link and the underlying economic forces are hard to pin down due to difficulty in measuring the degree of technological progress over time. Our goal is to fill this gap by constructing indices of technological progress at the aggregate and sectoral level that are consistently available over long periods of time.

Patent statistics are a useful starting point (Griliches, 1998). A major obstacle in inferring the degree of technological progress from patent data is that patents vary greatly in their technical and economic significance. While measures such as citations a patent receives in the future have been used to address this obstacle, these metrics have significant disadvantages. First, patent citations are only consistently recorded by the USPTO in patent documents after 1947. Prior to 1947, citations sometimes appear inside the text of the patent document, but they are much less common than in the post-war era.<sup>1</sup> Second, citations tend to take discrete values (the median post-1947 patent has 3 citations in a 10-year forward window). Third, citations rely on the discretion of the inventor or the patent examiner in choosing which prior patents to cite, or whether they are aware of the existence of closely related patents.<sup>2</sup>

Given these shortcomings, we instead propose a new indicator of patent importance that is similar in spirit and can be constructed by analyzing the text of patent documents. Our indicators require no other inputs besides the text of the patent document; hence they are consistently available for the entire history of US patents, spanning nearly two centuries of innovation (1840–2010). Further, unlike citations, there is limited discretion in how much of the patent document is written (besides the claims), since it pertains to the technical description of the innovation.

We start by leveraging natural language processing techniques to create links between each new invention and the set of existing and subsequent patents. Specifically, we construct measures of textual similarity to quantify commonality in the topical content of each pair of patents. We then identify an important patent as one whose content is distinct from prior

---

<sup>1</sup>For instance, consider patent 388,116 issued to William Seward Burroughs on August 1888 for a ‘calculating machine’, one of the precursors to the modern computer. Burroughs’ patent has just three citations as of March 2018. Similarly, patent 174,465 issued to Graham Bell for the telephone in February 1876 has the first recorded citation in 1956 (from patent 2,807,666). Until March 2018, it has received a total of 10 citations. These issues are not confined to the pre-1947 period: one of the first computer patents 2,668,661 issued in 1954 to George Stibitz at Bell Labs has just 15 citations as of March 2018.

<sup>2</sup>As an example, patent 6,368,227 for “Method of swinging on a swing”, issued to Steven Olson (aged 5) in April 2002, has 11 citations as of June 2018. It is cited, for example, by patent 8,420,782 for “Modular DNA-binding domains and methods of use”; patent 8,586,526 for “DNA-binding proteins and uses thereof”; and patent 8,697,853 for “TAL effector-mediated DNA modification”. Many of these citations were added by the patent examiner.

patents (is novel), but is similar to future patents (is impactful). These innovations represent distinct improvements in the technological frontier and become the new foundation upon which subsequent inventions are built. If citation data were objectively determined and consistently available, a breakthrough innovation would receive a large number of future citations.

Several tests confirm the validity of our measure of patent importance. First, we identify a set of major technological breakthroughs of the 19<sup>th</sup> and 20<sup>th</sup> century using the help of research assistants. Our indicators of patent significance perform quite well in identifying these major technological breakthroughs. Next, focusing on the post-1947 sample when citations data are available, we find our indicator is significantly correlated with patent citations. More importantly however, we find that our text-based patent indicators are significant predictors of future citations—indicating that they provide a more timely assessment of a patent’s quality than citation counts. Last, we relate our indicator to measures of private values. Though we view our indicators as more likely to be measuring the scientific value of a patent, prior work has documented a strong correlation between patent citations (which form the inspiration for our measure) and measures of market value (e.g. Hall et al., 2005; Kogan et al., 2017). We find that our quality indicator is significantly correlated with the Kogan et al. (2017) measure of a patent’s economic value.

Next, we construct time-series indices that describe the arrival intensity of breakthrough innovations—at the aggregate and sectoral level—by counting the number of patents each year whose importance is in the top decile of our importance measure (breakthrough patents). Our aggregate innovation index uncovers three major technological waves: the second Industrial Revolution (mid- to late 19<sup>th</sup> century), the 1920s and 1930s, and the post-1980 period. Inventions related to electricity were important in the late 19<sup>th</sup> and early 20<sup>th</sup> century. Innovations in agriculture played an important role in the beginning of the 20<sup>th</sup> century, while advances in genetically modified food have peaked in the last two decades. Chemical and petroleum-related innovations were particularly important in the 1920s and 1930s. Computers and electronic products have peaked since the early 1990s.

## I. Measuring Patent Similarity

Here, we discuss how to measure the similarity between pairs of patent documents; aggregate these similarities into a patent-level measure of importance; and construct a time-series index of breakthrough innovations.

## A. Data

We build our dataset from two sources: the USPTO and Google’s patent search engine. Our final dataset includes the full text of over nine million patents over the period 1840–2010. The Online Appendix provides additional details on our data collection process, as well as the conversion of unstructured patent text data into a numerical format suitable for statistical analysis.

### 1. Definition

A key consideration in devising a similarity metric for a pair of text documents is to appropriately weigh words by their importance. It is more informative if terms such as ‘electricity’ and ‘petroleum’ enter more prominently into the similarity calculation than common words like ‘process’ or ‘inventor.’ In textual analysis, a leading approach to overweighting terms that are most diagnostic of a document’s topical content is the “term-frequency-inverse-document-frequency” (*TFIDF*) transformation of word counts:

$$TFIDF_{pw} \equiv TF_{pw} \times IDF_w. \quad (1)$$

The first component of the weight, term frequency (TF), is defined as

$$TF_{pw} \equiv \frac{c_{pw}}{\sum_k c_{pk}}, \quad (2)$$

and describes the relative importance of term  $w$  for patent  $p$ . It counts how many times term  $w$  appears in patent  $p$  adjusted for the patent’s length. The second component is the inverse document frequency (IDF) of term  $w$ , which is defined as

$$IDF_w \equiv \log \left( \frac{\# \text{ documents in sample}}{\# \text{ documents that include term } w} \right). \quad (3)$$

*IDF* measures the informativeness of term  $w$  by under-weighting common words that appear in many documents, as these are less diagnostic of the content of any individual document.

The product of these two terms, *TFIDF*, describes the importance of a given word or phrase  $w$  in a given document  $p$ . Words that appear infrequently in a document tend to have low *TFIDF* scores (due to low *TF*), as do common words that appear in many documents (due to low *IDF*). A high value of *TFIDF*<sub>pw</sub> indicates that term  $w$  appears relatively frequently in document  $p$  but does not appear in most other documents, thus conveying that word  $w$  is especially representative of document  $p$ ’s semantic content. Younger and Kuhn (2016) use this method to measure the pairwise patent-to-patent similarity across a large subset of the USPTO patents.

For our purposes, this weighting scheme is not ideal, since we are interested in the novelty or impact of patent  $p$ 's text content given the history of innovation leading up to the development of  $p$ . Consider for example Nikola Tesla's famous 1888 patent (number 381,968) of an AC motor, which was among the first patents to use the phrase "alternating current," a phrase used with great frequency throughout the 20th century. Standard  $IDF$  would sharply de-emphasize this term in the  $TFIDF$  vector representing Tesla's patent because so many patents subsequently used this phrase so intensively.  $TFIDF$  would therefore give a misleading, and quite inverted, portrayal of the patent's importance.

We therefore develop a modified version of the traditional  $TFIDF$  measure. In place of (3), we instead construct a retrospective version of inverse document frequency. We define the "backward- $IDF$ " of term  $w$  for patent  $p$ , (denoted by  $BIDF_{wp}$ ) as the log frequency of documents containing  $w$  in any patent granted *prior* to patent  $p$ :

$$BIDF_{wp} = \log \left( \frac{\# \text{ patents prior to } p}{1 + \# \text{ documents prior to } p \text{ that include term } w} \right). \quad (4)$$

This frequency measure evolves as a term becomes more or less widely used over time, reflecting the history of invention up to, but not beyond, the new patent's arrival.

Continuing with the Tesla example discussed above, consider measuring the similarity between Tesla's AC motor patent, and patent 4,998,526 assigned in 1990 to General Motors Corporation for an "Alternating current ignition system." An important question emerges: what is the most sensible  $IDF$  to use when calculating  $TFIDF$  similarity of these two patents? One possibility is to use  $BIDF$  for the year 1888 in the  $TFIDF$  of Tesla's patent, and  $BIDF$  as of 1990 for GM's patent. However, over the 102 years between these two patents, "alternating current" appears in tens of thousands of other patents. Thus, the use of "alternating current" by GM would be greatly down-weighted with a 1990  $BIDF$  adjustment, and thus the co-occurrence of "alternating current" in these two patents would have a small contribution to the pair's similarity. Given our goal of quantifying the impact of patents on future technological innovations, we calculate pairwise similarity by applying to *both* patent counts the  $BIDF$  corresponding to the *earlier* of the two patents. Thus, to calculate the similarity between the patent pair in this Tesla/GM example, the term frequencies of both are normalized by the 1888 backward- $IDF$ .

In sum, we construct the similarity between the patent pair  $(i, j)$  as follows. First, for both patents we construct our modified-version of the  $TFIDF$  for each term  $w$  in patent  $i$  as

$$TFBIDF_{w,i,t} = TF_{w,i} \times BIDF_{w,t}, \quad t \equiv \min(i, j) \quad (5)$$

and likewise for patent  $j$ . These are arranged in a  $W$ -vector  $TFBIDF_{i,t}$  where  $W$  is the size

of the set union for terms in pair  $(i, j)$ . Next, each  $TFBIDF$  vector is normalized to have unit length,

$$V_{i,t} = \frac{TFBIDF_{i,t}}{\|TFBIDF_{i,t}\|}. \quad (6)$$

Finally, we calculate the cosine similarity between the two normalized vectors:

$$\rho_{i,j} = V_{i,t} \cdot V_{j,t}. \quad (7)$$

Because  $TFBIDF$  is non-negative,  $\rho_{i,j}$  lies in the interval  $[0,1]$ . Patents that use the exact same set of words in the same proportion will have similarity of one, while patents with no overlapping terms have similarity of zero.

## 2. Descriptive statistics

Panel A of Figure 1 plots the distribution of our similarity score across patent pairs that are 0–20 years apart. We see that patents tend to be highly dissimilar, with only a small fraction of pairs very closely related. The median similarity score across patent pairs is 7.8%, whereas the average similarity score is 10.2%. In the right tail, the 90<sup>th</sup> and 95<sup>th</sup> percentiles of similarity scores are 17.6% and 22.9%, respectively. In network terminology, the patent system’s connectivity is sparse.

Citation linkages provide external validation for assessing the text-based similarity measure  $\rho_{i,j}$ . Panel B examines whether patent pairs with high  $\rho_{i,j}$  are more likely to be linked by a citation. Indeed, we see that the likelihood that patent  $j$  cites the earlier patent  $i$  is monotonically increasing in the similarity  $\rho_{i,j}$  between the two patents.

## 3. Examples

Figure 2 provides a few examples of patents’ similarity networks. To simplify the presentation, and also illustrate the advantages of our method in the early parts of the sample, we focus on four patents from the 19th century. For each of these patents, the figure plots the set of prior and subsequent patents (filed within a period of five years) that have a cosine similarity of 50% or greater with the focal patent.

The patent at the top left part of the figure (US 4,750) is one of the first patents associated with the sewing machine, issued to 1846 to Elias Howe Jr. The patent is for the lockstitch, a manufacturing process still in use today. This patent is not significantly connected to any prior patents. By contrast, it is relatively closely related to sixteen patents, all for improvements in the sewing machine, that were filed over the next five years. Many of these subsequent patents were owned by either Elias Howe, or three companies, Wheeler & Wilson, Grover and Baker, and I. M. Singer, who together formed the first patent pool in American industry in

1856 (Lampe and Moser, 2010).

The patent on the top right (US 493,426) is one of the earliest patents associated with cinematography. The patent is issued to Thomas Edison, for exhibiting ‘photographs of moving objects’, by Thomas Edison, and is one of the earliest film projectors. The patent is highly similar to two prior patents and twelve subsequent patents, filed within five years. Most of the subsequent patents are related to cinematography. Among them, three are for a ‘kinetographic’ camera, one of the early precursors of the film camera.

The patent at the bottom, left part of the figure (US 161,739) is one of the early patents issued to Graham Bell and eventually led to the invention of the telephone. We can see that it is quite similar to four prior ‘telegraph’ patents filed over the previous five years. It is also related to eleven patents filed over the next five years, one of which is Graham Bell’s famous ‘telephone’ patent (174,465). Last, the patent on the bottom right is a random patent (US 222,189) for improvements in the cover of petroleum lamps. Within a five-year span, it is related to seven prior patents and five subsequent patents, all of which refer to improvements in lamps. In brief, our examples show that our similarity measure identifies meaningful connections between patents.

## II. Important Patents

Novel patents are those that are conceptually distinct from their predecessors, and therefore rely less on prior art. Impactful patents are those which influence future scientific advances, manifested as high similarity with subsequent innovations. The main idea in this paper is that an important patent is one that is *both novel and impactful*.

### A. Definition

We measure a patent’s novelty as its dissimilarity with the existing patent stock at the time it was filed. We start from a measure of ‘backward similarity’, defined as

$$BS_j^\tau = \sum_{i \in \mathcal{B}_{j,\tau}} \rho_{j,i}, \quad (8)$$

where  $\rho_{i,j}$  is the pairwise similarity of patents  $i$  and  $j$  defined in equation (7) and  $\mathcal{B}_{j,\tau}$  denotes the set of “prior” patents filed in the  $\tau$  calendar years prior to  $j$ ’s filing. Patents with low backward similarity are dissimilar to the existing patent stock. They deviate from the state of the art and are therefore novel. We consider a backward-looking window of  $\tau = 5$  years in our baseline importance measure—henceforth denoted by  $BS_j$ .

Next, we measure a patent’s impact by its “forward similarity,” defined as

$$FS_j^\tau = \sum_{i \in \mathcal{F}_{j,\tau}} \rho_{j,i}, \quad (9)$$

where  $\mathcal{F}_{j,\tau}$  denotes the set of patents filed over the next  $\tau$  calendar years following patent  $j$ ’s filing. The forward similarity measure in (9) estimates of the strength of association between the patent and future technological innovation over the next  $\tau$  years.

A patent might have high forward similarity because it changes the course of future innovation. Or, it might be part of scientific regime shift that was catalyzed by a predecessor patent. The “alternating current” example highlights this difference. Nikola Tesla’s patent has a high forward similarity because it dictated the course of future electronics, but was very different from any prior patents. The General Motors patent’s similarity with future AC-related patents merely reflects that it is part of a mainstream technology—it has a high similarity both backward and forward. Majorly important patents—those with a large influence on future technologies and that deviate from the status quo—are more likely to represent scientific breakthroughs.

The distinction between these two patents emerges when we compare forward versus backward similarity for a given patent. That is, our indicator of patent importance combines forward and backward similarity to identify patents that are both novel and impactful:

$$q_j^\tau = \frac{FS_j^\tau}{BS_j}. \quad (10)$$

This indicator attaches higher scientific value to patents that are both novel relative to their predecessors and are influential for subsequent research.

Our indicator of patent importance largely follows the logic behind indicators based on future citations. Specifically, the numerator in (10) is the total similarity with future patents—which is directly analogous to the sum of future citations. The denominator scales the forward similarity score by the novelty of the patent—since, presumably, patents should be citing the earliest relevant prior patents.<sup>3</sup>

The key advantage of our indicator is that it relies only on the text of the patent document, and is therefore broadly available. However, it also has limitations. Existing computational constraints limit the window over which we can compare patents; hence, our measure may under-weigh innovations whose impact took time to materialize. Further, our algorithm is reliant on the digitization quality of the patent document; patents with inaccurate text will be

---

<sup>3</sup>In contemporaneous work, Ashtor (2019) also constructs a measure of quality using textual similarity (estimated using latent semantic analysis). In contrast to this paper, Ashtor (2019) identifies high quality patents as those that have high similarity to contemporaneous and prior patents, uses only the claims portion of the document, and does not use any information on future similarity.



less similar to other (prior or subsequent) patents. Last, our algorithm in identifying impactful patents may be affected by shifting within-field propensity to patent (for instance, the recent rise of software patents). Such patents may appear to be impactful—as they are related to subsequent similar patents—and are more likely to be cited. Our measure shares this potential shortcoming with patent citations. One possibility is to extend our definition of impact (9) to include non-patent literature, which we leave to future work.

## B. Validation

Next, we conduct three validation checks for our importance measure. First, we identify a list of historically significant patents and examine how they score in terms of our importance indicators. Second, we relate our indicators to forward patent citations, a common measure of patent quality in the innovation literature. Last, we examine the correlation between our importance indicators and market values.

### 1. Historically significant patents

We compile a list of approximately 250 ‘historically significant patents’ based on online lists. For instance, the USPTO has a “Significant Historical Patents of the United States” list. Our list targets indisputable important and radical inventions of the last 200 years, beginning with the telegraph and internal combustion engine, and ending with stem cells, Google’s PageRank algorithm and gene transfer. The full list of patents and online sources is provided in Appendix Table A.1.

For each of these radical inventions we report their percentile rank in terms of our importance measure (10); for instance, a value of 0.90 indicates that the patent is in the top 10% of most important patents. We compute a patent’s rank using three approaches. First, we use the unconditional distribution. Second, we rank patents after subtracting the mean importance measure within each cohort (issue year). Removing cohort fixed effects helps eliminate factors that affect patents symmetrically, such as shifts in language or variation in the quality of the digitized patent documents. Second, we compute ranks within cohort. Though this comparison is not very useful in constructing a time-series index of technological change, it clarifies the extent to which these indicators are useful for purely cross-sectional comparisons.

Overall, these 250 patents score highly in terms of our importance indicator. The mean rank of these patents is 0.74 when using the unconditional distribution, or 0.78 when performing either adjustment. That said, our importance measure does miss some important inventions, such as pasteurization and Morse Code which our importance measure ranks at the bottom 20%.

One way to assess the performance of our measure is to compare how these patents rank in

terms of their citations. Since many of these patents are filed during the period when citation data is not broadly available, we extend the horizon for citations and measure them using the entire sample. Naturally this skews the comparison in favor of citations since they are measured over a significantly longer horizon than our indicators.

We find that our importance measure moderately outperforms citations: the average rank assigned to these important patents is 0.74, compared to 0.55 for citations when citations are measured using the full sample. The difference shrinks when these indicators are demeaned using year fixed effects, but is not fully eliminated—0.78 for importance versus 0.74 for citations. Removing time fixed effects leads to similar results as comparing patents within cohorts (mean rank 0.78 for importance versus 0.69 for citations). Appendix Figure A.1 summarizes these findings. We conclude that our text-based importance indicators are considerably more informative than patent citations in comparing patents across different cohorts, especially once we consider that our importance indicator can be computed using only 10 years of data—as opposed to several decades in many cases.

## 2. Patent citations

We next investigate the relation between our importance measure and patent citations, a commonly used metric of impact. We focus on patents issued after 1947, as this is the period when citations are consistently recorded by the USPTO. Panel A of Figure 3 illustrates the correlation using binned scatter plots; Appendix Table A.2 reports the corresponding regression estimates.

The first row of Panel A of Figure 3 reveals a strong positive contemporaneous correlation between patent importance and forward citations. Specifically, we first consider forward windows of  $\tau = 1, 5,$  and  $10$  years for both citations and importance. The correlation is consistently economically significant across horizons  $\tau$ . Comparing two patents in the same technology class that are issued to the same entity in the same year, we find that increasing the importance measure from the median to the 90<sup>th</sup> percentile results in 1.5 additional citations, relative to the median of 3 citations, when importance and citations are measured over the next 10 years after the patent application is filed.

One way to examine the additional information content of our measure relative to citation counts measured over the same horizon is to examine whether it predicts *future* patent citations. The second row of Panel A shows that this is indeed the case. We plot the predictive relation between our text-based quality measured in the  $0\text{-}\tau$  year window after filing, versus all citations in years  $\tau + 1$  and beyond. We control for the number of citations over the same period for which importance is measured. In all cases, we find an strong positive association between our near-term quality measure and long-term future citations. Comparing two patents in the same class, issued to the same entity in the same year, we see that an increase in the patent

importance from the median to the 90<sup>th</sup> percentile predicts 20-25% more future citations relative to the median. The results strongly suggest that our importance measure incorporates information faster than forward citations.

### 3. Estimates of market value

We next discuss the relation between patent importance and market valuations. Market values are by definition private values; they measure the present value of pecuniary benefits to the holder of the patent. By contrast, our importance measure is designed to ascertain the scientific importance of the patent. The relation between market value and scientific importance can be ambiguous. For instance, a patent may represent only a minor scientific advance while being very effective in restricting competition, thus generating large private rents (see, e.g. Abrams et al., 2013). With that caveat in mind, we next examine the relation between our importance measure and the estimate of patent value of Kogan et al. (2017)—henceforth KPSS. The KPSS measure,  $\hat{V}_j$ , infers the value of patent  $j$  (in dollars) from stock market reaction to the patent grant. KPSS interpret this measure as an ex-ante measure of the private value of the patent.

Panel B of Figure 3 graphically illustrates this correlation; Appendix Table A.3 reports the corresponding regression estimates. Patent importance is positively and statistically significantly correlated to the KPSS estimate of market value. Focusing on two patents in the same class that are issued to the same firm in the same year, increasing the importance measure from the median to the 90<sup>th</sup> percentile results in 0.23–0.47% increase in patent values. Though these estimates may appear relatively modest, they are comparable in magnitude to the relation between patent values and forward citations (see, e.g. Kogan et al., 2017). Further, Appendix Table A.3 shows that the correlation remains significant once we include as additional controls the number of forward citations the patent receives over the same horizon that importance is measured—which supports the conclusion that our measure incorporates additional information to patent citations.

## III. Indices of Technological Progress

Our goal here is to construct indices of breakthrough innovations—that is, innovations that are in the right tail of our importance measure—that span the USPTO sample (1840–2010).

### A. Construction

One challenge is that time-series fluctuations in (10) are mechanically affected by factors such as shifts in language; the fact that the retrospective document frequency measure (4) is changing over time so terms become less novel over time; and the fact that the number of patents is

rapidly expanding over time. Given that these issues likely affect most patents symmetrically, we adjust (10) by removing patent cohort issue year fixed effects. After this adjustment, we define a ‘breakthrough’ patent as one that falls in the top 10% of the unconditional distribution of importance estimated as the ratio of 10-year forward to 5-year backward similarity. We then construct a time series index as the number of breakthrough inventions granted in each year, divided by US population. The implicit assumption in our methodology is that shifts in language are likely to symmetrically affect all patents and will thus be absorbed by the fixed effect—which is mainly identified by the non-breakthrough patents. We also construct indices of innovation at the sector level using the probabilistic mapping between patent technology classifications (CPC) and industry classifications constructed by Goldschlag et al. (2016).

To validate our methodology, we show that our technology indices are significantly related to measured productivity, both at the aggregate as well as sectoral level. As we discuss in Section F in the Appendix, a one-standard deviation increase in our index is associated with 0.5% to 2% higher annual productivity growth over the next ten years. Similarly, sectors that have breakthrough innovations experience faster growth in productivity than sectors that do not: a one-standard deviation increase in our innovation index is associated with 1% higher annual productivity growth over the next five years.

## B. Aggregate Index

Panel A of Figure 4 plots the resulting time-series of breakthroughs per capita. Our index identifies three main innovation waves, lasting from 1870 to 1880; 1920 to 1935; and from 1985 to the present. The first peak corresponds to the beginning of the second industrial revolution, which saw technological advances such as the telephone and electric lighting and improvements in railroads. The second peak corresponds to advances in manufacturing, particularly in plastics and chemicals, consistent with the evidence of Field (2003). The latest wave of technological progress includes revolutions in computing, genetics, and telecommunication.

Constructing an innovation index has proven challenging in the past. In one approach, Shea (1999) constructs an index of per-capita patent counts, which is plotted in Panel B. Patents per capita is essentially flat from 1870–1930, dips from 1930–1980, and displays a significant spike post-1980. There are reasons to be skeptical that such an index indeed measures the degree of underlying progress, since it implicitly assumes that all patents are equally valuable. One common adjustment to simple patent counts is to weigh patents by their forward citations. Panel C (black line) plots the resulting time-series when our index methodology is instead constructed from 10-year forward citations. Due to the limitations of citation data, this series essentially identifies no innovation prior to 1940s. Only when citations are measured over the entire sample (blue line) does the index take non-zero values in the pre-WW2 period, but

even then the levels dwarf the values of the index post-1980. Given that the importance of inventions in the 1850–1940 era are at least comparable to the those in the last two decades (see, e.g. Gordon, 2016), this pattern mostly reflects the limitations of forward citations as a measure of patent importance.

Kogan et al. (2017) construct a time-series index that is based on the estimated market values of patents that are granted. Their index is plotted in Panel D. Their index has the advantage that it provides a dollar estimate of the value of innovation output in a given year. However, it is confined to the universe of publicly traded firms, thereby omitting innovations by private firms, non-profit institutions and the government. Moreover, it is not available prior to 1927, since information on stock prices is readily available only after this year.

## C. Sectoral Indices

Figure 5 plots time-series indices of industry innovation at the 3-digit NAICS level. We see that the origin of breakthrough patents has varied considerably over time. In the 1840–70 period, we see that the most important inventions took place in engineering and construction, consumer goods, and manufacturing. An example of an important invention in construction according to our importance measure is the ‘Bollman Bridge’ (patent 8,624)—the first successful all-metal bridge design widely used for railroads. Other important advances in this period occur in textiles: examples include various versions of sewing and knitting machines (patents 7,931; 7,296; 7,509; and 60,310).

Starting around 1870, many more patents that score high in terms of our measure are related to electricity, with some of the most important patents relating to the production of electric light (203,844; 210,380; 215,733; 210,213; 200,545; 218,167). The same period saw the invention of a revolutionary method of communication: the telephone. It is comforting that most of the patents associated with the telephone are among the top 1%.<sup>4</sup>

Another industry that accounted for a significant share of important patents during the 1860-1910 period is transportation. Many of the patents that fall in the top 1% include improvements in railroads (207,538; 218,693; 422,976; and 619,320), and their electrification (178,216; 344,962; 403,969; 465,407). Most importantly, the turn of the century saw the invention of the airplane. In addition to the Wright brothers’ original patent (821,393), several other airplane patents also score highly in terms of our importance indicator (1,107,231; 1,279,127; 1,307,133; 1,307,134) as well as patents related to air balloons and the Zeppelin (678,114 and 864,672). Innovations in construction methods continue to play a role in this period, such as those that are related to the use of concrete (618,956; 647,904; 764,302; 654,683; 747,652; and 672,176) in the construction of buildings, roads and pavements.

---

<sup>4</sup>Patents: 161,739; 174,465; 178,399; 186,787; 201,488; 213,090; 220,791; 228,507; 230,168; 238,833; 474,230; 203,016; 222,390.

In the first half of the 20th century, chemistry emerges as a major generator of important patents, many describing inventions of plastic compounds. Among our breakthrough inventions is the patent for bakelite (942,699), the world’s first fully synthetic plastic. This innovation opened the floodgates to a torrent of now-familiar synthetic plastics, including the invention in the 1930’s of PVC by Waldo Semon (patents 1,929,453 and 2,188,396) and nylon by Wallace H. Carothers (patent 2,071,250), all of which score in the top 5%. Other important chemistry patents in the 1950s include drugs, such as Nystatin (2,797,183); improvements in the production of penicillin (2,442,141 and 2,443,989); Enovid, the first oral contraceptive (2,691,028); and Tetracycline, one of the most prescribed broad spectrum antibiotics (2,699,054).

The 1950s are marked by the harnessing of nuclear energy for civilian purposes. Enrico Fermi’s patents on the development of the nuclear reactor all score highly.<sup>5</sup> Subsequent to the 1950’s, a large fraction of the important patents identified by our measure are in the area of Instruments and Electronics, and are related to the arrival of the Information Age. One of the most important patents according to our measure is the invention of the first microchip by Robert Noyce in 1961 (patent 2,981,877). During the 1970s, firms such as IBM, Xerox, Honeywell, AT&T, and Sperry Rand develop some of the major innovations in computing. Xerox, for example, produced for several high-scoring inventions such as patent 4,558,413 for a management system software; patent 4,899,136 for improvements in computer user interface; patent 4,437,122 for bitmap graphics; and patents 3,838,260 and 3,938,097 for improvements in the interface between computer memory and the processor.

In the 1980s and 1990s, several important patents that pertain to computer networks emerge among the set of breakthrough patents.<sup>6</sup> Improvements in genetics comprise a significant fraction of the most important patents in the 1980–2000 period. A few early examples that fall in the top 1% of the unconditional distribution according to our importance indicator are: patent 4,237,224 for recombinant DNA methods; patents 4,683,202; 4,683,195, and 4,965,188 for the PCR method for rapidly copying DNA segments with high fidelity and at low cost; patent 4,736,866 for genetically modified animals; and patent 4,889,818 for heat-stable DNA-replication enzymes.

## IV. Conclusion

We use textual analysis of high-dimensional data from patent documents to create a new measure of technological innovation that allows us to characterize the evolution of technological waves over the entire 1840–2010 period across a broad set of sectors.

---

<sup>5</sup>Patents 2,206,634; 2,836,554; 2,524,379; 2,852,461; 2,708,656; 2,768,134; 2,780,595; 2,798,847; 2,807,581; 2,807,727; 2,813,070; 2,837,477; 2,931,762.

<sup>6</sup>Patents 4,800,488; 4,823,338; 4,827,411; 4,887,204; 5,249,290; 5,341,477; 5,544,322; and 5,586,260.

## References

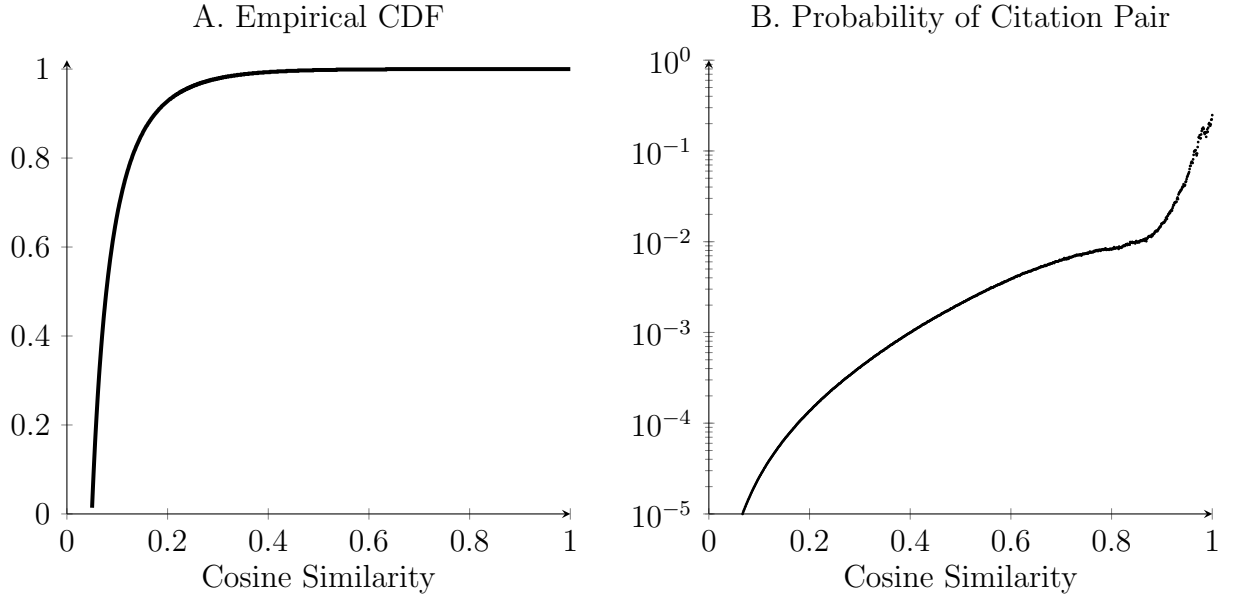
- Abrams, D. S., U. Akcigit, and J. Popadak (2013). Patent value and citations: Creative destruction or strategic disruption? Working Paper 19647, National Bureau of Economic Research.
- Ashtor, J. H. (2019). Investigating cohort similarity as an ex ante alternative to patent forward citations. *Journal of Empirical Legal Studies* 16(4), 848–880.
- Basu, S., J. G. Fernald, and M. S. Kimball (2006). Are technology improvements contractionary? *American Economic Review* 96(5), 1418–1448.
- Berkes, E. (2016). Comprehensive universe of u.s. patents (cusp): Data and facts. Working paper, Northwestern University.
- Field, A. J. (2003). The most technologically progressive decade of the century. *American Economic Review* 93(4), 1399–1413.
- Goldschlag, N., T. J. Lybbert, and N. J. Zolas (2016). An ‘algorithmic links with probabilities’ crosswalk for uspc and cpc patent classifications with an application towards industrial technology composition. CES Discussion Paper 16-15, U.S. Census Bureau.
- Gordon, R. (2016). *The Rise and Fall of American Growth: The U.S. Standard of Living since the Civil War*. The Princeton Economic History of the Western World. Princeton University Press.
- Griliches, Z. (1998, January). *Patent Statistics as Economic Indicators: A Survey*, pp. 287–343. University of Chicago Press.
- Hall, B. H., A. B. Jaffe, and M. Trajtenberg (2005). Market value and patent citations. *The RAND Journal of Economics* 36(1), pp. 16–38.
- Jorda, O. (2005, March). Estimation and inference of impulse responses by local projections. *American Economic Review* 95(1), 161–182.
- Kendrick, J. W. (1961). *Productivity Trends in the United States*. National Bureau of Economic Research, Inc.
- Kogan, L., D. Papanikolaou, A. Seru, and N. Stoffman (2017). Technological innovation, resource allocation, and growth\*. *The Quarterly Journal of Economics* 132(2), 665–712.
- Lampe, R. and P. Moser (2010). Do patent pools encourage innovation? evidence from the nineteenth-century sewing machine industry. *The Journal of Economic History* 70(4), 898–920.

Shea, J. (1999). What do technology shocks do? In *NBER Macroeconomics Annual 1998*, volume 13, NBER Chapters, pp. 275–322. National Bureau of Economic Research, Inc.

Younge, K. and J. M. Kuhn (2016). Patent-to-patent similarity: A vector space model. Working paper.



Figure 1: Pairwise similarity and citation linkages



Panel A plots the empirical CDF of our similarity measure  $\rho_{i,j}$  across patent citation pairs. Panel B plots the conditional probability that patent  $i$  cites an earlier patent  $j$  as a function of the text-based similarity score between the two patents,  $\rho_{i,j}$ , computed in equation (7) in the main text. Specifically, we bin patent pairs  $(i, j)$  in terms of their cosine similarity, and then compute the average propensity of a citation link—that is, we estimate  $E[\mathbf{1}_{i,j}|\rho_{i,j}]$ , where  $\mathbf{1}_{i,j}$  is a dummy variable that takes the value one if patent  $j$  cites patent  $i$  (where patent  $i$  is filed prior to patent  $j$ ). For computational reasons, we exclude similarity pairs with  $\rho_{i,j} \leq 5\%$ . Figure uses data only post 1945, since citations were not consistently recorded prior to that year.

Figure 2: Examples of Similarity Networks

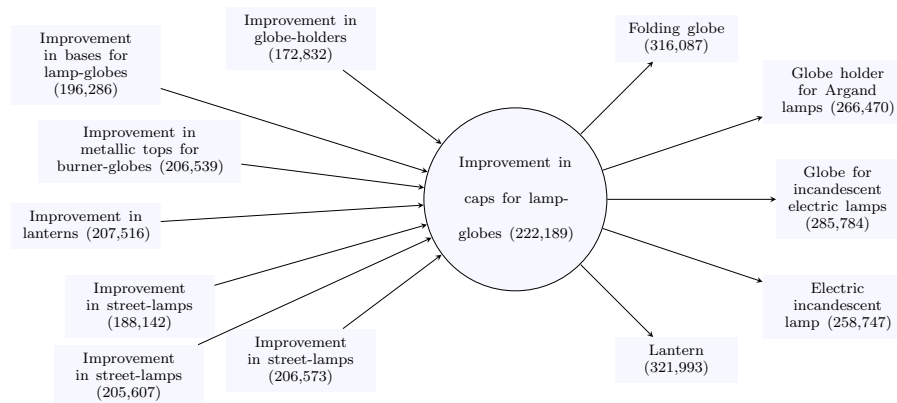
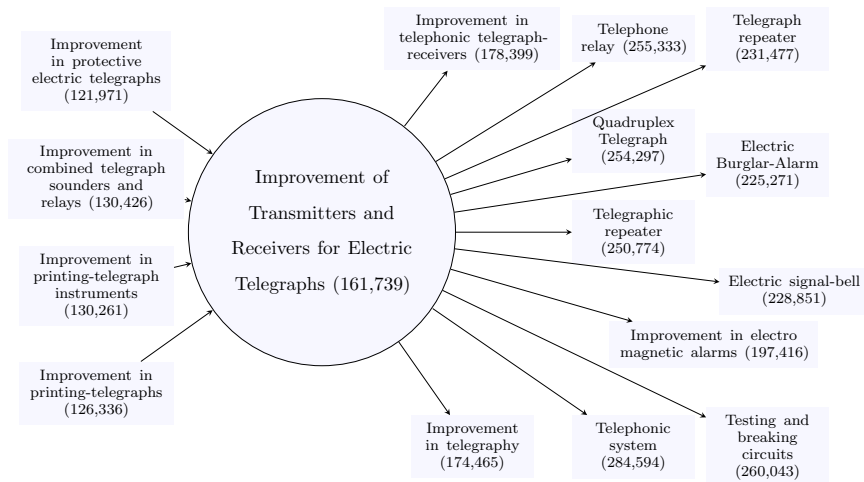
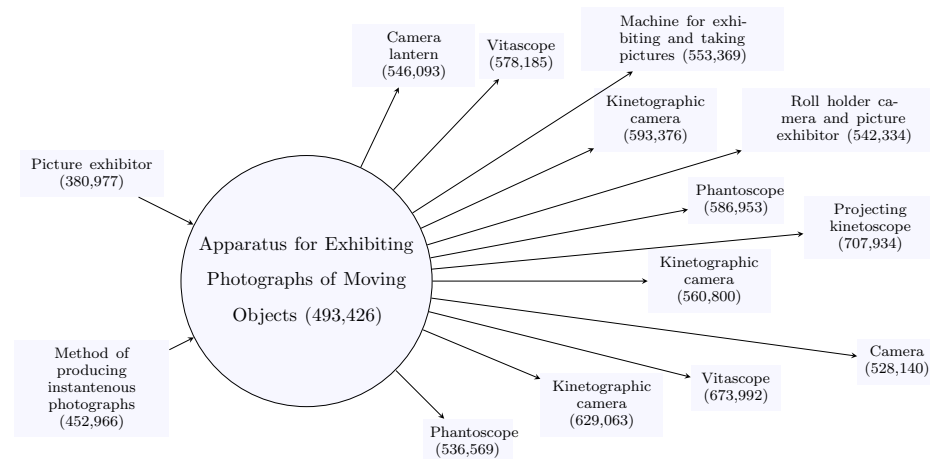
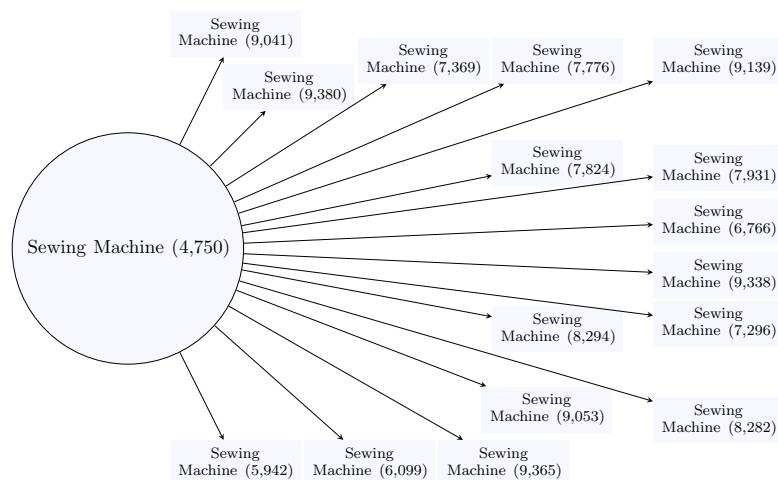
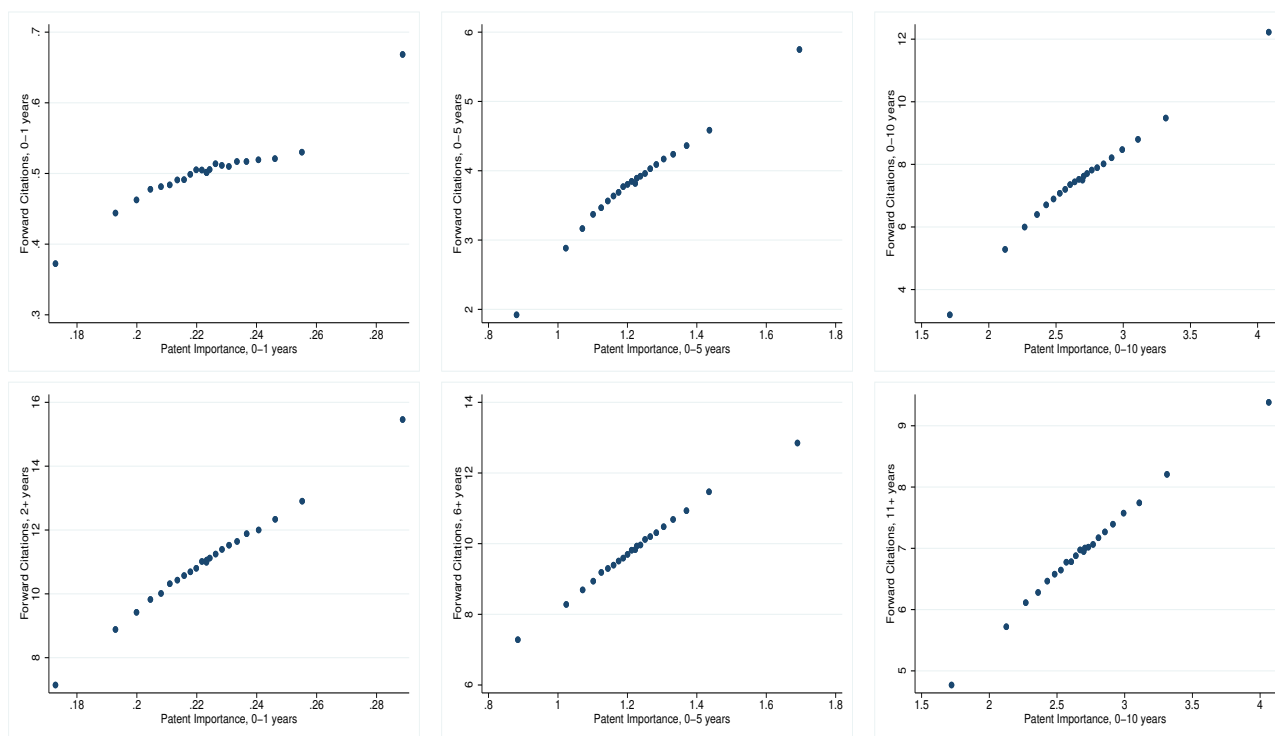


Figure displays the similarity network for four patents: the patent for the first sewing machine (top left); one of the earlier patents for moving pictures (top right); one of the early patents that led to the telephone (bottom left) and a randomly chosen patent from the 1800s (bottom right). In plotting the similarity links, we restrict attention to patents pairs filed at most five years apart and with a cosine similarity greater than 50%.

**Figure 3: Patent Importance: Validation**

A. Patent Citations



B. Patent Value (KPSS)

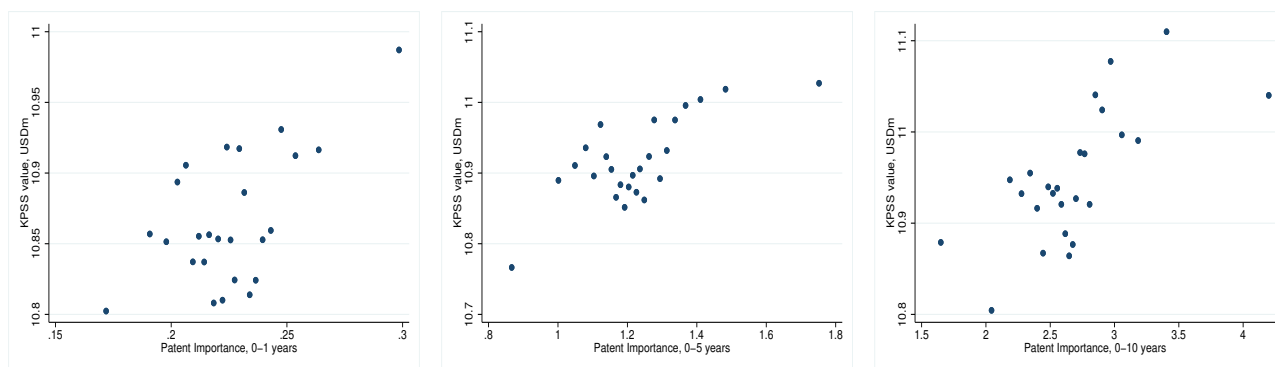
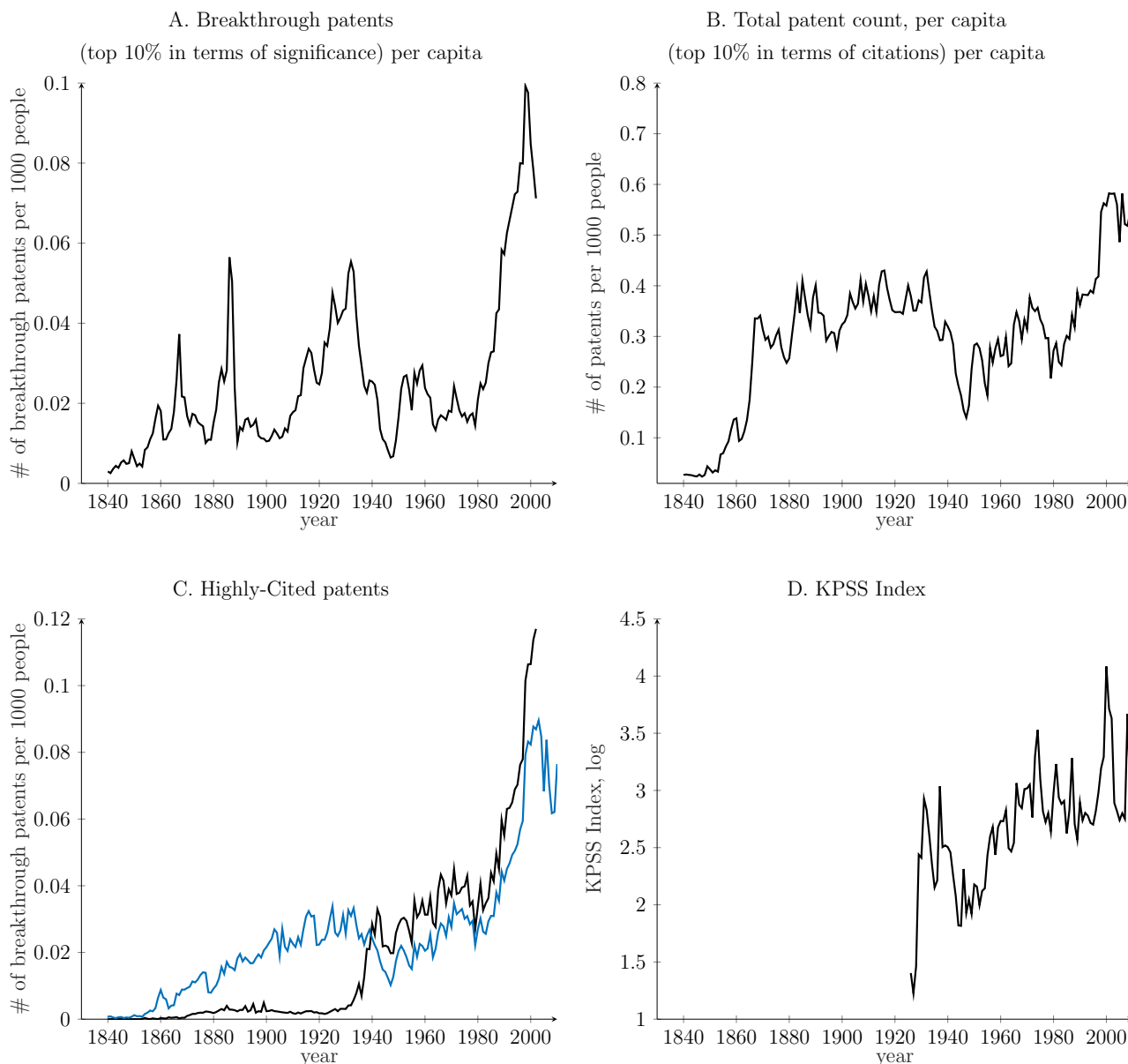


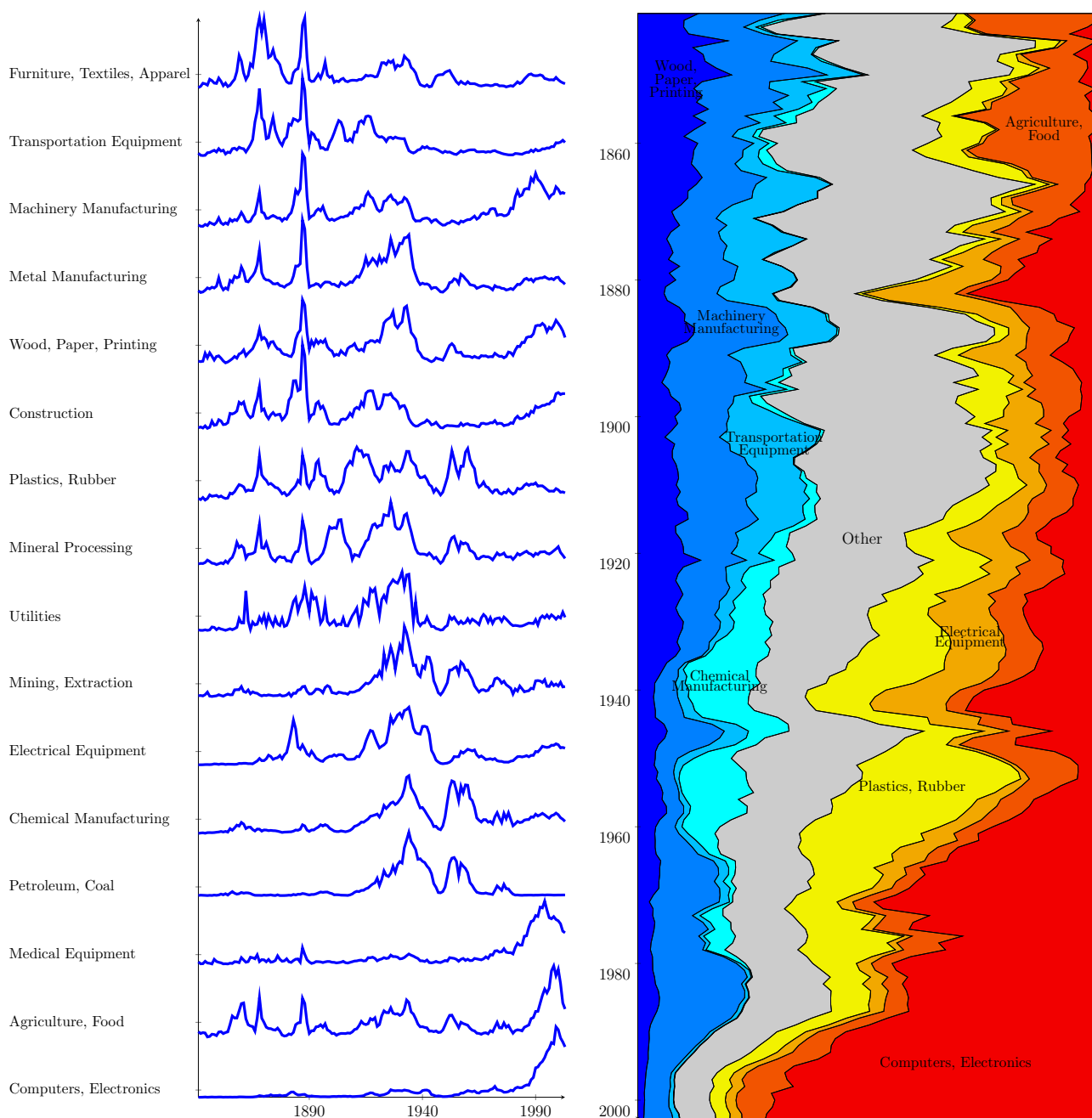
Figure plots the relation between our importance measure (the ratio of forward to backward patent similarity) to the number of forward citations (Panel A) or the Kogan et al. (2017) estimate of patent value (Panel B). In the first row of Panel A, both the patent importance measure and forward citations are measured over the same horizon. In the second row of Panel B, we plot the predictive relation between our importance measure and future citations; for these specifications, we also control for the number of citation the patent has received over the same horizon that our importance measure is computed. To construct the figure, we group observations into 25 bins (cutoff at every other percentile of the quality distribution). Within each bin, we average citation and text-based importance measures after controlling for technology class and assignee-by-grant year fixed effects. See Appendix Tables A.2 and A.3 for the corresponding regression tables.

**Figure 4: Technological Innovation over the Long Run: New vs Existing Indicators**



Panel A plots the number of breakthrough patents per capita. Breakthrough patents are those that fall in the top 10% of the unconditional distribution of our importance measure, where importance is defined as the ratio of the 10 year forward to the 5 year backward similarity, net of year fixed effects. Panel B plots the total number of patents, scaled by population. In Panel C we plot, in black (blue), the number of patents that fall in the top 10% of the unconditional distribution of forward citations—measured over the next 10 years (entire sample), net of year fixed effects—again scaled by US population. while Panel D plots the KPSS Index (the sum of the estimated market value of patents scaled by the total capitalization of the stock market).

Figure 5: Breakthrough Innovation Across Industries



We plot the per capita number of breakthrough patents across industries. Industries are defined based on NAICS codes. Breakthrough patents are those that fall in the top 10% of our baseline importance measure (defined as the ratio of the 10-yr forward to the 5-yr backward similarity) net of issue year fixed effects. We construct industry indices using the CPC4 to NAICS crosswalk constructed by Goldschlag et al. (2016).

# Appendix

We briefly overview our conversion of unstructured patent text data into a numerical format suitable for statistical analysis. To begin, we build our collection of patent documents from two sources. The first is the USPTO patent search website, which records all patents beginning from 1976. Our web crawler collected the text content of patents from this site, which includes patent numbers 3,930,271 through 9,113,586. The records in this sample are comparatively easy to process as they are available in HTML format with standardized fields.

For patents granted prior to 1976, we collect patent text from our second main datasource, Google’s patent search engine. For the pre-1976 patent records, we recover all of the fields listed above with the exception of inventor/assignee addresses (Google only provides their names), examiner, and attorney. Some parts of our analysis rely on firm-level aggregation of patent assignments. We match patents to firms by firm name and patent assignee name. Our procedure broadly follows that of Kogan et al. (2017) with adaptations for our more extensive sample. In addition to the citation data we scrape from Google, we obtain complementary information on patent citations from Berkes (2016) and the USPTO. The data in Berkes (2016) includes citations that are listed inside the patent document and which are sometimes missed by Google. Nevertheless, the likelihood of a citation being recorded is significantly higher in the post-1945 than in the pre-1945. When this consideration is relevant, we examine results separately for the pre- and post-1945 periods.

To represent patent text as numerical data, we convert it into a *document term matrix* (DTM), denoted  $C$ . Columns of  $C$  correspond to words and rows correspond patents. Each element of  $C$ , denoted  $c_{pw}$ , counts the number of times a given one-word phrase (indexed by  $w$ ) is used in a particular patent (indexed by  $p$ ), after imposing a number of filters to remove stop words, punctuation, and so forth. We provide a detailed step-by-step account of our DTM construction in Appendix IV. Our final dictionary includes 1,685,416 terms in the full sample of over nine million patents.

The next section provides additional details on the data construction, including the process through which we convert the text of patent documents to a format that is amenable to constructing similarity measures.

## A. Text Data Collection, Additional Details

The Patent Act of 1836 established the official US Patent Office and is the grant year of patent number one.<sup>7</sup> We construct a dataset of textual content of US patent granted during the 180 year period from 1836-2015. Our dataset is built on two sources.

---

<sup>7</sup>The first patent was granted in the US in 1790, but of the patents granted prior to the 1836 Act, all but 2,845 were destroyed by fire.

The first is the USPTO patent search website. This site provides records for all patents beginning in 1976. We designed a web crawler collect the text content of patents over this period, which includes patent numbers 3,930,271 through 9,113,586. We capture the following fields from each record:

- |                        |                        |                        |
|------------------------|------------------------|------------------------|
| 1. Patent number (WKU) | 7. Assignee addresses  | 13. Backward citations |
| 2. Application date    | 8. Family ID           | 14. Examiner           |
| 3. Granted date        | 9. Application number  | 15. Attorney           |
| 4. Inventors           | 10. US patent class    | 16. Abstract           |
| 5. Inventor addresses  | 11. CPC patent class   | 17. Claims             |
| 6. Assignees           | 12. Intl. patent class | 18. Description        |

The only information available from USPTO that we do not store are image files for a patent’s “figure drawing” exhibits.

For patents granted prior to 1976, the USPTO also provides bulk downloads of .txt files for each patent. The quality of this data is inferior to that provided by the web search interface in three ways. First, the text data is recovered from image files of the original patent documents using OCR scans. OCR scans often contain errors. These generally arise from imperfections in the original images that lead to errors in the OCR’s translation from image to text. Going backward in time from 1976, the quality of OCR scans deteriorates rapidly due to lower quality typesetting. Second, the bulk download files do not use a standardized format which makes it difficult to parse out the fields listed above.

Rather than using the USPTO bulk files, we collect text of pre-1976 patents from our second main datasource, Google’s patent search engine. Like post-1976 patents from USPTO, Google provides patent records in an easy-to-parse HTML format that we collect with our web crawler. Furthermore, inspection of Google records versus 1) OCR files from the USPTO and 2) pdf images of patents that are the source of the OCR scans, reveals that in this earlier period Google’s patent text is more accurate than the OCR text in USPTO bulk data. From Google’s pre-1976 patent records, we recover all of the fields listed above with the exception of inventor/assignee addresses (Google only provides their names), examiner, and attorney.

## **B. Cleaning Post-1976 USPTO Data**

Next, we conduct a battery of checks to correct data errors. For the most part, we are able to capture and parse of patent text from the USPTO web interface without error. When there are errors, it is almost always the case that the patent record was incompletely captured, and this occurs for one of two reasons. The first reason is that the network connection was interrupted during the capture and the second is that the patent record on the UPSTO website

is itself incomplete (in comparison with PDF image files of the original document, which are also available from USPTO via bulk download).

Our primary data cleaning task was to find and complete any partially captured patent records. First, we find the list of patent numbers (WKUs) that are entirely missing from our database, and re-run our capture program until all have been recovered. Many of the missing records that we find are explicitly labeled as “WITHDRAWN” at the USPTO.<sup>8</sup> Next, we identify WKUs with an entirely missing value for the abstract, claims, or description field. Fortunately, we find this to be very infrequent, occurring in less than one patent in 100,000, making it easy for us to correct this manually.

Next, a team of research assistants (RA’s) manually checked 3,000 utility patent records, 1,000 design patent records, and 1,000 plant patents records against their PDF image files. The RA task is to identify any records with missing or erroneous information in the reference, abstract, claims, or description fields. To do this, they manually read the original pdf image for the patent and our digitally captured record. We identify patterns in partial text omission and update our scraping algorithm to reflect these. We then re-ran the capture program on all patents and confirmed that omissions from the previous iteration were corrected.

### C. Cleaning Pre-1976 Google Data

Fortunately, we find no instances of missing WKU’s or incomplete text from Google web records. Next, we assess the accuracy of Google’s OCR scans by manually re-scanning a random sample of 1,000 pre-1976 patents using more recent (and thus more accurate) ABBYY OCR software than was used for most of Google’s image scans. We compare the ABBYY scan to the pdf image to confirm the scan content is complete, the compare the frequency of garbled terms in our scan versus that OCR text from Google. The distribution of pairwise cosine similarities in our ABBYY text and Google’s OCR is reported below.

---

<sup>8</sup>Withdrawn information can be found at <https://www.uspto.gov/patents-application-process/patent-search/withdrawn-patent-numbers>.



Cosine Similarity	
mean	0.957
std	0.073
P1	0.701
P5	0.863
P10	0.900
P25	0.951
P50	0.977
P75	0.991
P90	0.996
P95	0.998
P99	0.999
N	1000

Only 10% of sampled Google OCR records have a correlation with ABBYY below 90%.

Next, we manually compare both our OCR scans and those from Google against the pdf image. We find that garble rate for ABBYY OCRed is 0.025 on average, with standard deviation of 0.029. We find that Google has only slightly more frequent garbling than our ABBYY scans. Of the term discrepancies in the two sets of scans, around 52% of these correspond to a garbled ABBYY records and 83% to a garbled Google record. We ultimately conclude that Google’s OCR error frequency is acceptable for use in our analysis.

## D. Conversion from Textual to Numeric Data

We convert the text content of patents into numerical data for statistical analysis. To do this, we use the NLTK Python Toolkit to parse the “abstract,” “claims,” and “description” sections of each patent into individual terms. We strip out all non-word text elements, such as punctuation, numbers, and HTML tags, and convert all capitalized characters to lowercase. Next, we remove all occurrences of 947 “stop words,” which include prepositions, pronouns, and other words that carry little semantic content.<sup>9</sup>

<sup>9</sup>We construct our stop word list as the union of terms in the following commonly used lists:

<http://www.ranks.nl/stopwords>  
<https://dev.mysql.com/doc/refman/5.1/en/fulltext-stopwords.html>  
<https://code.google.com/p/stop-words/>  
<http://www.lextek.com/manuals/onix/stopwords1.html>  
<http://www.lextek.com/manuals/onix/stopwords2.html>  
<http://www.webconfs.com/stop-words.php>  
<http://www.text-analytics101.com/2014/10/all-about-stop-words-for-text-mining.html>  
[http://www.nlm.nih.gov/bsd/disted/pubmedtutorial/020\\_170.html](http://www.nlm.nih.gov/bsd/disted/pubmedtutorial/020_170.html)  
<https://pypi.python.org/pypi/stop-words>  
<https://msdn.microsoft.com/zh-cn/library/bb164590>  
<http://www.nltk.org/book/ch02.html> (NLTK list)

The remaining list of “unstemmed” (that is, without removing suffixes) unigrams amounts to a dictionary of 35,640,250 unique terms. As discussed in Gentzkow, Kelly, and Taddy (2017), an important preliminary step to improve signal-to-noise ratios in textual analysis is to reduce the dictionary by filtering out terms that occur extremely frequently or extremely infrequently. The most frequently used words show up in so many patents that they are uninformative for discriminating between patent technologies. On the other hand, words that show up in only a few patents can only negligibly contribute to understanding broad technology patterns, while their inclusion increases the computational cost of analysis.<sup>10</sup>

We apply filters to retain influential terms while keeping the computational burden of our analysis at a manageable level, and focus on the number of distinct patents and calendar years in which terms occur. A well known attribute of text count data is its sparsity—most terms show up very infrequently—and the table shows that this pattern is evident in patent text as well. We exclude terms that appear in fewer than twenty out of the more than nine million patents in our sample. These eliminate 33,954,834 terms, resulting in a final dictionary of 1,685,416 terms.<sup>11</sup>

After this dictionary reduction, the entire corpus of patent text is reduced in a  $D \times W$  numerical matrix of term counts denoted  $C$ . Matrix row  $d$  corresponds to patent (WKU)  $d$ . Matrix column  $w$  corresponds the  $w^{th}$  term in the dictionary. Each matrix element  $c_{dw}$  the count of term  $w$  in patent  $d$ .

Pairwise similarities constitute a high-dimensional matrix of approximate dimension 9 million  $\times$  9 million. To reduce the computational burden, we set similarities below 5% to zero. This affects 93.4% of patent pairs. Patents with such low text similarity are, for all intents and purposes, completely unrelated, yet introduce a large computational load in the types of analyses we pursue. Replacing these approximate zeros with similarity scores of exactly zero achieves large computational gains by allowing us to work with sparse matrix representations that require substantially less memory. Our empirical findings are insensitive to this threshold as they are driven primarily by the highest similarity pairs. In experiments with similarity cutoffs ranging from 5% to 10%, we find results that are quantitatively indistinguishable.

## E. Matching Patents to Firms

Much of our analysis relies on firm-level aggregation of patent assignments. We match patents to firms by merging firm names and patent assignee names. Our procedure broadly follows

---

<sup>10</sup>Filtering out infrequent words also removes garbled terms, misspellings, and other errors, as their irregularity leads them to occur only sporadically.

<sup>11</sup>The table also shows that there are some terms that appear in almost all patents. Examples of the most frequently occurring words (that are not in the stop word lists) are “located,” “process,” and “material.” Because these show up in most patents they are unlikely to be informative for statistical analysis. These terms are de-emphasized in our analysis through the *TFIDF* transformation.

that of Kogan et al. (2017) with adaptations for our more extensive sample.

The first step is extracting assignee names from patent records. For post-1976 data we use information from the USPTO web search to identify assignee names. Due to the high data quality in this sample, assignee extraction is straightforward and highly accurate. For pre-1976, we use assignee information from Google patent search. While it is easy to locate the assignee name field thanks to the HTML format, Google’s assignee names are occasionally garbled by the OCR.

Next, we clean the set of extracted assignee names. There are 766,673 distinct assignees in patents granted since 1836. Most of the assignees are firm names and those that are not firms are typically the names of inventors. We clean assignee name garbling using fuzzy matching algorithms. For example, the assignee “international business machines” also appears as an assignee under the names “ininternational business machines,” “international businesss machines,” and “international business machiness.” Garbled names are not uncommon, appearing for firms as large as GE, Microsoft, Ford Motor, and 3M.

We primarily rely on Levenshtein edit distance between assignees to identify and correct erroneous names. There are two major challenges to overcome in name cleaning. The first choosing a distance threshold for determining whether names are the same. As an example, the assignees “international business machines” (recorded in 103,544) and “ibm” (recorded in 547 patents) have a large Levenshtein distance. To address cases like this, we manually check the roughly 3,000 assignee names that have been assigned at least 200 patents, correcting those that are variations on the same firm name (including the IBM, GE, Microsoft, Ford, and 3M examples). Next, for each firm on the list of most frequent assignees, we calculate the Levenshtein distance between this assignee name and the remaining 730,000+ assignee names, and manually correct erroneous names identified by the list of assignees with short Levenshtein distances.

The second challenge is handling cases in which a firm subsidiary appears as assignee. For example, the General Motors subsidiary “gm global technology operations” is assigned 8,394 patents. To address this, we manually match subsidiary names from the list of top 3,000+ assignees to their parent company by manually searching Bloomberg, Wikipedia, and firms’ websites.

After these two cleaning steps, and after removing patents with the inventor as assignee, we arrive at 3,036,859 patents whose assignee is associated with a public firm in CRSP/Compustat, for a total of 7,467 distinct cleaned assignee firm names. We standardized these names by removing suffixes such as “com,” “corp,” and “inc,” and merge these with CRSP company names. Again we manually check the merge for the top 3,000+ assignees, and check that name changes are appropriately addressed in our CRSP merging step. Finally, we also merge our patent data with Kogan et al. (2017) patent valuation data for patents granted between 1926

and 2012.

## F. Breakthrough Innovation and Measured Productivity

Here, we relate our innovation indices to measured productivity.

### 1. Aggregate Productivity

For the post-war sample, we use the aggregate TFP measure constructed by Basu et al. (2006), which is available over the 1948-2018 period. For the earlier sample, we measure productivity using output per hour data collected by Kendrick (1961), which is available for the 1889 to 1957 period. Following Jorda (2005), we estimate:

$$\frac{1}{\tau} (x_{t+\tau} - x_t) = a_0 + a_\tau \log \text{BreakthroughIndex}_t + c_\tau \mathbf{Z}_t + u_{t+\tau}, \quad (11)$$

where  $x_t$  is log productivity,  $\text{BreakthroughIndex}_t$  refers to our innovation index, and  $\mathbf{Z}_t$  is a vector of controls that includes the log number of patents per capita and the level of productivity. We consider horizons of up to  $\tau = 10$  years and adjust the standard errors for serial correlation using the Newey-West procedure with  $\tau + 1$  lags. All independent variables are normalized to unit standard deviation. To ensure that we are not capturing pre-existing trends, we also examine negative values of  $\tau$ .

Panel A of Figure A.3 presents the results of estimating (11) for the post-war sample. We see that a one-standard deviation increase in our index is associated with 0.5 percent faster annual TFP growth, with some delay. This is substantial given that the standard deviation in measured TFP growth over this period is 1.8%. Panel B shows the results for the earlier sample. Again, we see that a one-standard deviation increase in our innovation index is associated with an increase in labor productivity growth of approximately 1.5–2% per year—compared to an annual standard deviation of 5.2% for labor productivity growth.

### 2. Sector-level Productivity

First, we examine how the distribution across technology class of breakthrough patents varies over time. Panel A of Appendix Figure A.2 shows the technology classes in which breakthrough inventions originated has varied substantially over the last 170 years. By contrast, Panel B shows that the composition of technology classes among all patents has remained relatively stable over time.

We next construct indices of innovation at the sector level. One issue that arises is how to map patents to industries in a way that is independent of the presence of an explicit assignee, since clean assignee identity and names are notoriously difficult to pin down. To address

this, we exploit the mapping between patent technology classifications (CPC) and industry classifications constructed by Goldschlag et al. (2016). Because this is a probabilistic mapping (there is no one-to-one correspondence between CPC and industry codes), we assign a fraction of each patent to industry codes based on the given probability weights associated with its (4-digit) CPC technology classification.<sup>12</sup>

Figure 5 plots time-series indices of industry innovation at the 3-digit NAICS level. We see that the origin of breakthrough patents has varied considerably over time, consistent with our prior results. Inventions related to electricity were important in the late 19<sup>th</sup> and early 20<sup>th</sup> century. Innovations in agriculture played an important role in the beginning of the 20<sup>th</sup> century, while advances in genetically modified food have peaked in the last two decades. Chemical and petroleum-related innovations were particularly important in the 1920s and 1930s. Computers and electronic products have peaked since the early 1990s. We next examine whether our industry indices are related to measured sectoral productivity.

Panel A of Figure A.4 presents our results for the period from 1987 to the present. We use estimates of multi-factor productivity at the NAICS 4-digit level from the Bureau of Labor Statistics (BLS), which covers 86 manufacturing industries. We then estimate a panel analogue of equation (11),

$$\frac{1}{\tau}(x_{i,t+\tau} - x_{i,t}) = a_0 + a_\tau \log \text{BreakthroughIndex}_{i,t} + c_\tau \mathbf{Z}_{i,t} + u_{i,t+\tau}, \quad (12)$$

except that now  $\mathbf{Z}_{i,t}$  also includes time and industry fixed effects. Standard errors are clustered by industry. Given the shorter length of this sample, we consider horizons of  $\tau = 1 \dots 5$  years.

We find a strongly statistically positive relation between our innovation index and future productivity growth—while the relation with past productivity growth is insignificant. In terms of magnitudes, a one-standard deviation increase in our innovation index is associated with approximately 1–1.2% higher productivity growth per year, over the next 5 years.

Panel B performs a similar exercise for the earlier sample. We use the labor productivity data collected by Kendrick (1961), which covers 62 manufacturing industries for the years 1899, 1909, 1919, 1937, 1947, and 1954. Since the data is only available at discrete periods, we modify our approach: for each period  $(t, t + \tau)$ , we regress the annualized difference in log labor productivity on the log of the accumulated level of innovation (number of breakthrough patents) in  $t \pm 2$  years. We again see a strong and statistically significant relation between our industry innovation indices and measured productivity: a one standard deviation increase in our innovation index is associated with a 1.4% higher growth rate in measured productivity

---

<sup>12</sup>Two caveats are in order. First, this mapping is based on post-1970 data, whereas our analysis spans the entire period since the 1840s. Hence, there might be measurement error in our index since we assign a fraction of patents to each of the industries that map to a CPC classification based on the weights estimated from only part of the sample. Second, this mapping is primarily available for manufacturing industries—which are however the industries that patent most heavily.

over the next period.

For comparison, we also construct a corresponding index based on citations (measured over a 10 year horizon). Examining Panels A and B of Appendix Figure A.5, we see that there is no statistically significant relation between the citations-based index and industry productivity in either sample period.

# Appendix Material

Table A.1: Important Patents

Patent	Year	Inventor	Invention	Citations	Percentile Ranks						Source
					No Adjustment			Remove year FE			
					Quality	Citations		Quality	Citations		
					(total)	(0-10)	(0-10)	(total)	(0-10)	(0-10)	
1647	1840	Samuel F. B. Morse	Morse Code	2	0.03	-	0.29	0.03	0.64	0.81	Reference
3237	1843	Nobert Rillieux	Sugar Refining	0	0.80	-	-	0.84	0.64	0.44	Reference
3316	1843	Samuel F. B. Morse	telegraphy wire	0	0.97	-	-	0.99	0.64	0.44	Reference
3633	1844	Charles Goodyear	Vulcanized Rubber	3	0.99	-	0.38	0.98	0.64	0.88	Reference
4453	1846	Samuel F. B. Morse	telegraph battery	0	1.00	-	-	0.99	0.64	0.44	Reference
4750	1846	Elias Howe, Jr.	Sewing Machine	1	1.00	-	0.17	0.99	0.64	0.70	Reference
4834	1846	Benjamin Franklin Palmer	Artificial Limb	0	0.99	-	-	0.87	0.64	0.44	Reference
4848	1846	Charles T. Jackson	Anesthesia	0	0.98	-	-	0.75	0.64	0.44	Reference
4874	1846	Christian Frederick Schonbein	Guncotton	0	0.97	-	-	0.69	0.64	0.44	Reference
5199	1847	Richard M. Hoe	Rotary Printing Press	0	0.99	-	-	0.80	0.64	0.42	Reference
5711	1848	M. Waldo Hanchett	Dental Chair	1	1.00	-	0.17	0.99	0.64	0.70	Reference
5942	1848	John Bradshaw	Sewing Machine	0	1.00	-	-	0.98	0.64	0.44	Reference
6099	1849	Morey/Johnson	Sewing Machine	1	1.00	-	0.17	0.99	0.64	0.69	Reference
6281	1849	Walter Hunt	Safety Pin	0	1.00	-	-	0.94	0.64	0.42	Reference
6439	1849	John Bachelder	Sewing Machine	0	1.00	-	-	0.97	0.64	0.42	Reference
7296	1850	D.M. Smith	Sewing Machine	0	1.00	-	-	1.00	0.64	0.40	Reference
7509	1850	J. Hollen	Sewing Machine	0	1.00	-	-	1.00	0.64	0.40	Reference
7931	1851	Grover and Baker	Sewing Machine	0	1.00	-	-	0.99	0.64	0.40	Reference
8080	1851	John Gorrie	Ice Machine	0	0.99	-	-	0.35	0.64	0.40	Reference
8294	1851	Isaac Singer	Sewing Machine	0	1.00	-	-	0.98	0.64	0.40	Reference
9300	1852	Lorenzo L. Langstroth	Beehive	1	1.00	-	0.17	0.85	0.64	0.69	Reference
13661	1855	Isaac M. Singer	Shuttle Sewing Machine	1	0.95	-	0.17	0.03	0.50	0.63	Reference
15553	1856	Gail Borden, Jr.	Condensed Milk	0	0.99	-	-	0.92	0.63	0.34	Reference
17628	1857	William Kelly	Iron and Steel Manufacturing	0	0.99	-	-	0.85	0.54	0.35	Reference
18653	1857	H.N. Wadsworth	Toothbrush	6	0.98	-	0.58	0.61	0.54	0.94	Reference
23536	1859	Martha Coston	System of Pyrotechnic Night Signals	1	0.97	-	0.17	0.60	0.64	0.58	Reference
26196	1859	James J. Mapes	Artificial Fertilizer	1	0.99	-	0.17	0.94	0.64	0.58	Reference
31128	1861	Elisha Graves Otis	Elevator	1	0.98	-	0.17	0.78	0.41	0.46	Reference
31278	1861	Linus Yale, Jr.	Lock	10	0.96	-	0.72	0.60	0.41	0.94	Reference
31310	1861	Samuel Goodale	Moving Picture Peep Show Machine	0	0.99	-	-	0.95	0.41	0.18	Reference



Table A.1: Important Patents (cont)

Patent	Year	Inventor	Invention	Citations	Percentile Ranks						Source
					No Adjustment			Remove year FE			
					Quality	Citations		Quality	Citations		
					(total)	(0-10)	(0-10)	(total)	(0-10)	(0-10)	
36836	1862	Richard J. Gatling	Machine Gun	3	0.95	0.22	0.38	0.21	0.83	0.82	Reference
43465	1864	Sarah Mather	Submarine Telescope	0	0.94	-	-	0.03	0.41	0.40	Reference
46454	1865	John Deere	Plow	0	0.99	-	-	0.60	0.43	0.41	Reference
53561	1866	Milton Bradley	Board Game	2	1.00	-	0.29	1.00	0.51	0.81	Reference
59915	1866	Pierre Lallement	Bicycle	0	0.99	-	-	0.86	0.51	0.41	Reference
78317	1868	Alfred Nobel	Dynamite	4	0.65	-	0.46	0.09	0.64	0.92	Reference
79265	1868	C. Latham Sholes	Typewriter	1	0.93	-	0.17	0.81	0.64	0.69	Reference
79965	1868	Alvin J. Fellows	Spring Tape Measure	2	0.82	-	0.29	0.36	0.64	0.82	Reference
88929	1869	George Westinghouse	Air Brake	1	0.84	-	0.17	0.79	0.64	0.69	Reference
91145	1869	Ives W. McGaffey	Vacuum Cleaner	4	0.74	0.22	0.46	0.58	0.85	0.92	Reference
110971	1871	Andrew Smith Hallidie	Cable Car	1	0.80	0.22	0.17	0.79	0.83	0.67	Reference
113448	1871	Mary Potts	Sad Iron	3	0.67	-	0.38	0.55	0.41	0.87	Reference
127360	1872	J.P. Cooley, S. Noble	Toothpick-making machine	0	0.75	-	-	0.68	0.39	0.39	Reference
129843	1872	Elijah McCoy	Improvements in Lubricators for Steam-Engines	1	0.73	-	0.17	0.64	0.39	0.66	Reference
135245	1873	Louis Pasteur	Pasteurization	0	0.54	-	-	0.26	0.37	0.38	Reference
141072	1873	Louis Pasteur	Manufacture of Beer and Treatment of Yeast	1	0.48	-	0.17	0.18	0.37	0.66	Reference
157124	1874	Joseph F. Glidden	Barbed Wire	1	0.94	-	0.17	0.95	0.38	0.65	Reference
161739	1875	Alexander Graham Bell	Telephone	7	0.99	-	0.62	0.99	0.38	0.96	Reference
171121	1875	George Green	Dental Drill	2	0.97	0.22	0.29	0.98	0.83	0.79	Reference
174465	1876	Alexander Graham Bell	Telephone	6	1.00	0.37	0.58	1.00	0.92	0.95	Reference
178216	1876	Alexander Graham Bell	Telephone	0	1.00	-	-	1.00	0.39	0.38	Reference
178399	1876	Alexander Graham Bell	Telephone	2	0.99	0.22	0.29	0.99	0.83	0.79	Reference
186787	1877	Alexander Graham Bell	Electric Telegraphy	0	1.00	-	-	1.00	0.37	0.37	Reference
188292	1877	Chester Greenwood	Earmuffs	17	0.92	-	0.84	0.91	0.37	0.99	Reference
194047	1877	Nicolaus August Otto	Internal Combustion Engine	1	0.73	-	0.17	0.48	0.37	0.65	Reference
200521	1878	Thomas Alva Edison	Phonograph	12	0.91	0.37	0.77	0.84	0.91	0.98	Reference
201488	1878	Alexander Graham Bell	Telephone	2	1.00	-	0.29	1.00	0.36	0.78	Reference
203016	1878	Thomas Alva Edison	Speaking Telephone	15	1.00	0.37	0.82	1.00	0.91	0.99	Reference
206112	1878	Thaddeus Hyatt	Reinforced Concrete	0	0.79	-	-	0.47	0.36	0.36	Reference
220925	1879	Margaret Knight	Paper-Bag Machine	4	0.88	0.49	0.46	0.64	0.95	0.90	Reference

Table A.1: Important Patents (cont)

Patent	Year	Inventor	Invention	Citations	Percentile Ranks						Source
					No Adjustment			Remove year FE			
					Quality	Citations	Quality	Citations	Quality	Citations	
					(total)	(0-10)	(0-10)	(total)	(0-10)	(0-10)	
222390	1879	Thomas Alva Edison	Improvement in carbon telephones	16	1.00	-	0.83	1.00	0.36	0.99	Reference
223898	1880	Thomas Alva Edison	First Incandescent Light	20	0.99	-	0.87	0.99	0.41	0.99	Reference
224573	1880	Emile Berliner	Microphone	0	0.95	-	-	0.89	0.41	0.36	Reference
228507	1880	Alexander Graham Bell	Electric Telephone	3	1.00	0.37	0.38	1.00	0.92	0.85	Reference
237664	1881	Frederic E. Ives	Halftone Printing Plate	1	0.90	0.22	0.17	0.72	0.83	0.64	Reference
304272	1884	Ottmar Mergenthaler	Linotype	0	0.90	-	-	0.93	0.40	0.35	Reference
312085	1885	Edward J. Claghorn	Seat Belt	13	0.34	-	0.79	0.30	0.38	0.98	Reference
322177	1885	Sarah Goode	Folding Cabinet Bed	3	0.53	-	0.38	0.60	0.38	0.84	Reference
347140	1886	Elihu Thomson	Electric Welder	16	0.58	0.88	0.83	0.58	1.00	0.99	Reference
349983	1886	Gottlieb Daimler	Four Stroke Combustion Engine	4	0.98	-	0.46	0.99	0.39	0.89	Reference
371496	1887	Dorr E. Felt	Adding Machine	6	0.71	0.57	0.58	0.73	0.97	0.94	Reference
372786	1887	Emile Berliner	Phonograph Record	4	0.73	0.49	0.46	0.75	0.95	0.89	Reference
373064	1887	Carl Gassner, Jr.	Dry Cell Battery	3	0.28	-	0.38	0.11	0.38	0.84	Reference
382280	1888	Nikola Tesla	A. C. Induction Motor	2	0.80	0.22	0.29	0.89	0.83	0.76	Reference
386289	1888	Miriam Benjamin	Gong and Signal Chair for Hotels	0	0.50	-	-	0.53	0.41	0.34	Reference
388116	1888	William S. Burroughs	Calculator	3	0.76	-	0.38	0.85	0.41	0.84	Reference
388850	1888	George Eastman	Roll Film Camera	1	0.85	-	0.17	0.93	0.41	0.62	Reference
395782	1889	Herman Hollerith	Computer	1	0.54	0.22	0.17	0.66	0.83	0.61	Reference
400665	1889	Charles M. Hall	Aluminum Manufacture	2	0.85	0.22	0.29	0.94	0.83	0.76	Reference
415072	1889	Starley/Owen	Tandem Bicycle	1	0.63	-	0.17	0.77	0.42	0.61	Reference
430212	1890	Hiram Stevens Maxim	Smokeless Gunpowder	0	0.54	-	-	0.73	0.46	0.34	Reference
430804	1890	Herman Hollerith	Electric Adding Machine	2	0.80	0.22	0.29	0.93	0.84	0.76	Reference
447918	1891	Almon B. Strowger	Telephone Exchange	81	0.56	-	0.98	0.81	0.48	1.00	Reference
453550	1891	John Boyd Dunlop	Pneumatic Tyres	1	0.78	0.22	0.17	0.94	0.84	0.61	Reference
468226	1892	William Painter	Bottle Cap	7	0.73	-	0.62	0.94	0.34	0.94	Reference
472692	1892	G.C. Blickensderfer	Typewriting Machine	4	0.28	0.37	0.46	0.58	0.91	0.88	Reference
492767	1893	Edward G. Acheson	Carborundum	12	0.24	-	0.77	0.53	0.44	0.98	Reference
493426	1893	Thomas Alva Edison	Motion Picture	1	0.77	-	0.17	0.95	0.44	0.60	Reference
504038	1893	Whitcomb L. Judson	Zipper	6	0.24	-	0.58	0.53	0.44	0.93	Reference
536569	1895	Charles Jenkins	Phantoscope	0	0.87	-	-	0.96	0.34	0.31	Reference

Table A.1: Important Patents (cont)

Patent	Year	Inventor	Invention	Citations	Percentile Ranks						Source
					No Adjustment			Remove year FE			
					Quality	Citations		Quality	Citations		
					(total)	(0-10)	(0-10)	(total)	(0-10)	(0-10)	
549160	1895	George B. Selden	Automobile	0	0.69	-	-	0.87	0.34	0.31	Reference
558393	1896	John Harvey Kellogg	Cereal	3	0.60	-	0.38	0.67	0.49	0.83	Reference
558719	1896	C.B. Brooks	Street Sweeper	2	0.61	0.37	0.29	0.68	0.92	0.75	Reference
558936	1896	Joseph S. Duncan	Addressograph	3	0.33	0.22	0.38	0.25	0.84	0.83	Reference
586193	1897	Guglielmo Marconi	Radio	4	0.83	0.57	0.46	0.87	0.97	0.88	Reference
589168	1897	Thomas A. Edison	Motion Picture Camera	0	0.63	-	-	0.61	0.49	0.31	Reference
608845	1898	Rudolf Diesel	Diesel Engine	8	0.77	-	0.66	0.76	0.47	0.95	Reference
621195	1899	Ferdinand Graf Zeppelin	Dirigible	1	0.72	-	0.17	0.52	0.35	0.57	Reference
644077	1900	Felix Hoffmann	Aspirin	1	0.71	-	0.17	0.41	0.46	0.58	Reference
661619	1900	Valdemar Poulsen	Magnetic Tape Recorder	15	0.84	0.69	0.82	0.74	0.98	0.98	Reference
708553	1902	John P. Holland	Submarine	1	0.75	-	0.17	0.54	0.45	0.57	Reference
743801	1903	ÉMary Anderson	Windscreen Wiper	2	0.35	-	0.29	0.07	0.51	0.73	Reference
745157	1903	Clyde J. Coleman	Electric Starter	1	0.91	-	0.17	0.91	0.51	0.57	Reference
764166	1904	Albert Gonzales	Railroad Switch	0	0.67	-	-	0.59	0.52	0.30	Reference
766768	1904	Michael J. Owens	Glass Bottle Manufacturing	7	0.76	0.64	0.62	0.74	0.98	0.94	Reference
775134	1904	KC Gillette	Razor (with removable blades)	4	0.92	0.49	0.46	0.95	0.95	0.87	Reference
808897	1906	Willis H. Carrier	Air Conditioning	21	0.66	0.22	0.88	0.72	0.84	0.99	Reference
815350	1906	John Holland	Submarine	0	0.71	-	-	0.78	0.54	0.28	Reference
821393	1906	Orville Wright	Airplane	19	1.00	0.22	0.86	1.00	0.84	0.99	Reference
841387	1907	Lee De Forest	Triode Vacuum Tube	5	0.29	0.22	0.52	0.23	0.85	0.90	Reference
921963	1909	Leonard H. Dyer	Automobile Vehicle	0	0.59	-	-	0.77	0.57	0.26	Reference
942809	1909	Leo H. Baekeland	Bakelite	3	0.89	0.22	0.38	0.97	0.84	0.80	Reference
970616	1910	Thomas A Edison	helicopter (never flown)	2	0.91	-	0.29	0.98	0.61	0.71	Reference
971501	1910	Fritz Haber	Ammonia Production	1	0.97	0.22	0.17	0.99	0.85	0.54	Reference
1000000	1911	Francis Holton	Non-Puncturable Vehicle Tire	2	0.83	-	0.29	0.93	0.60	0.71	Reference
1005186	1911	Henry Ford	Automotive Transmission	3	0.59	-	0.38	0.76	0.60	0.80	Reference
1008577	1911	Ernst F. W. Alexanderson	High Frequency Generator	6	0.50	0.69	0.58	0.65	0.99	0.92	Reference
1030178	1912	Peter Cooper Hewitt	Mercury Vapor Lamp	1	0.85	-	0.17	0.95	0.55	0.54	Reference
1082933	1913	William D. Coolidge	Tungsten Filament Light Bulb	28	0.73	0.57	0.92	0.90	0.97	0.99	Reference
1102653	1914	Robert H. Goddard	Rocket	58	0.42	0.49	0.97	0.62	0.95	1.00	Reference

Table A.1: Important Patents (cont)

Patent	Year	Inventor	Invention	Citations	Percentile Ranks						Source
					No Adjustment			Remove year FE			
					Quality	Citations		Quality	Citations		
					(total)	(0-10)	(0-10)	(total)	(0-10)	(0-10)	
1103503	1914	Robert Goddard	Rocket Apparatus	29	0.36	0.49	0.92	0.53	0.95	0.99	Reference
1113149	1914	Edwin H. Armstrong	Wireless Receiver	11	0.87	0.22	0.75	0.97	0.85	0.97	Reference
1115674	1914	Mary P. Jacob	Brassiere	1	0.53	-	0.17	0.76	0.59	0.53	Reference
1180159	1916	Irving Langmuir	Gas Filled Electric Lamp	13	0.78	0.64	0.79	0.94	0.98	0.97	Reference
1203495	1916	William D. Coolidge	X-Ray Tube	11	0.77	0.49	0.75	0.93	0.95	0.96	Reference
1211092	1917	William Coolidge	X-Ray Tube	7	0.94	0.22	0.62	0.99	0.84	0.92	Reference
1228388	1917	Frederick C Bargar	Fire Extinguisher	2	0.51	-	0.29	0.78	0.53	0.68	Reference
1254811	1918	Charles F. Kettering	Engine Ignition	1	0.50	-	0.17	0.78	0.60	0.51	Reference
1279471	1918	Elmer A. Sperry	Gyroscopic Compass	9	0.94	0.22	0.69	0.99	0.85	0.95	Reference
1360168	1920	Ernst Alexanderson	Antenna	4	0.92	-	0.46	0.98	0.62	0.83	Reference
1394450	1921	Charles P Strite	Bread Toaster	2	0.60	-	0.29	0.85	0.62	0.66	Reference
1413121	1922	John Arthur Johnson	Adjustable Wrench	0	0.05	-	-	0.06	0.63	0.20	Reference
1420609	1922	Glenn H. Curtiss	Hydroplane	2	0.68	-	0.29	0.89	0.63	0.65	Reference
1573846	1926	Thomas Midgley, Jr.	Ethyl Gasoline	3	0.36	0.22	0.38	0.78	0.84	0.72	Reference
1682366	1928	Charles F. Brannock	Foot Measuring Device	4	0.10	0.22	0.46	0.38	0.84	0.78	Reference
1699270	1929	John Logie Baird	Television / TV	11	0.55	-	0.75	0.91	0.48	0.94	Reference
1773079	1930	Clarence Birdseye	Frozen Food	10	0.53	0.22	0.72	0.92	0.84	0.93	Reference
1773080	1930	Clarence Birdseye	Frozen Food	18	0.60	-	0.86	0.94	0.45	0.97	Reference
1773980	1930	Philo T. Farnsworth	Television	29	0.83	0.91	0.92	0.98	1.00	0.99	Reference
1800156	1931	Erik Rotheim	Aerosol Spray Can	30	0.66	0.22	0.93	0.96	0.84	0.99	Reference
1821525	1931	Nielsen Emanuel	Hair Dryer	11	0.14	-	0.75	0.63	0.44	0.93	Reference
1835031	1931	Herman Affel	Coaxial cable	15	0.52	0.64	0.82	0.93	0.98	0.96	Reference
1848389	1932	Igor Sikorsky	Helicopter	5	0.41	-	0.52	0.91	0.42	0.78	Reference
1867377	1932	Otto F Rohwedder	Bread-Slicing Machine	2	0.09	-	0.29	0.57	0.42	0.52	Reference
1925554	1933	John Logie Baird	Color Television	1	0.38	-	0.17	0.90	0.37	0.33	Reference
1929453	1933	Waldo Semon	Rubber	56	0.83	0.96	0.97	0.99	1.00	1.00	Reference
1941066	1933	Edwin H. Armstrong	FM Radio	0	0.55	-	-	0.95	0.37	0.10	Reference
1948384	1934	Ernest O. Lawrence	Cyclotron	96	0.39	0.22	0.99	0.91	0.82	1.00	Reference
1949446	1934	William Burroughs	Adding and Listing Machine	1	0.09	0.22	0.17	0.62	0.82	0.31	Reference
1980972	1934	Lyndon Frederick	Krokodil	1	0.66	-	0.17	0.97	0.36	0.31	Reference

Table A.1: Important Patents (cont)

Patent	Year	Inventor	Invention	Citations	Percentile Ranks						Source
					No Adjustment			Remove year FE			
					Quality	Citations		Quality	Citations		
					(total)	(0-10)	(0-10)	(total)	(0-10)	(0-10)	
2021907	1935	Vladimir K. Zworykin	Television	18	0.56	0.37	0.86	0.95	0.91	0.95	Reference
2059884	1936	Leopold D. Mannes	Color Film	15	0.26	0.57	0.82	0.80	0.96	0.93	Reference
2071250	1937	Wallace H. Carothers	Nylon	231	0.79	0.97	1.00	0.98	1.00	1.00	Reference
2087683	1937	PT Farnsworth	Image Dissector	1	0.58	-	0.17	0.93	0.27	0.23	Reference
2153729	1939	Ernest H. Volwiler	Pentothal (General Anesthetic)	2	0.66	-	0.29	0.94	0.21	0.38	Reference
2188396	1940	Waldo Semon	Rubber	59	0.91	0.64	0.97	0.99	0.94	0.99	Reference
2206634	1940	Enrico Fermi	Radioactive Isotopes	99	0.78	0.99	0.99	0.97	1.00	1.00	Reference
2230654	1941	Roy J. Plunkett	TEFLON	49	0.48	0.93	0.96	0.91	0.99	0.99	Reference
2258841	1941	Jozsef Bir— Laszlo	Fountain Pen	20	0.05	0.84	0.87	0.31	0.98	0.94	Reference
2292387	1942	Markey/Antheil	Secret Communication System	71	0.33	0.37	0.98	0.86	0.74	0.99	Reference
2297691	1942	Chester F. Carlson	Xerography	738	0.11	0.91	1.00	0.58	0.99	1.00	Reference
2329074	1943	Paul Muller	DDT - Insecticide	48	0.15	0.97	0.96	0.68	1.00	0.98	Reference
2390636	1945	Ladislo Biro	Ball Point Pen	27	0.31	0.94	0.92	0.71	0.99	0.95	Reference
2404334	1946	Frank Whittle	Jet Engine	35	0.17	0.94	0.94	0.31	0.99	0.97	Reference
2436265	1948	Allen Du Mont	Cathode Ray Tube	18	0.54	0.88	0.86	0.64	0.98	0.91	Reference
2451804	1948	Donald L. Campbell	Fluid Catalytic Cracking	9	0.63	0.77	0.69	0.76	0.94	0.77	Reference
2495429	1950	Percy Spencer	Microwave	15	0.25	0.80	0.82	0.20	0.96	0.89	Reference
2524035	1950	John Bardeen	Transistor	132	0.79	1.00	0.99	0.90	1.00	1.00	Reference
2543181	1951	Edwin H. Land	Instant Photography	116	0.61	0.99	0.99	0.76	1.00	1.00	Reference
2569347	1951	William Shockley	Junction Transistor	140	0.64	1.00	0.99	0.79	1.00	1.00	Reference
2642679	1953	Frank Zamboni	Resurfacing Machine	16	0.41	0.57	0.83	0.53	0.85	0.89	Reference
2668661	1954	George R. Stibitz	Modern Digital Computer	14	0.96	0.77	0.80	0.99	0.94	0.86	Reference
2682050	1954	Andrew Alford	Radio Navigation System	3	0.68	-	0.38	0.81	0.09	0.39	Reference
2682235	1954	Richard Buckminster Fuller	Geodesic Dome	86	0.57	0.96	0.99	0.69	1.00	0.99	Reference
2691028	1954	Frank B. Colton	First Oral Contraceptive	4	0.90	-	0.46	0.97	0.09	0.48	Reference
2699054	1955	Lloyd H. Conover	Tetracycline	38	0.91	0.95	0.95	0.96	0.99	0.97	Reference
2708656	1955	Enrico Fermi	Atomic Reactor	196	0.96	1.00	1.00	0.99	1.00	1.00	Reference
2708722	1955	An Wang	Magnetic Core Memory	76	0.82	0.99	0.98	0.90	1.00	0.99	Reference
2717437	1955	George De Mestral	Velcro	258	0.39	0.98	1.00	0.35	1.00	1.00	Reference
2724711	1955	Gertrude Elion	Leukemia-fighting drug 6-mercaptopurine	1	0.78	0.22	0.17	0.88	0.26	0.13	Reference

Table A.1: Important Patents (cont)

Patent	Year	Inventor	Invention	Citations	Percentile Ranks						Source
					No Adjustment			Remove year FE			
					Quality	Citations	Quality	Citations	Quality	Citations	
					(total)	(0-10)	(0-10)	(total)	(0-10)	(0-10)	
2752339	1956	Percy L. Julian	Preparation of Cortisone	11	0.87	0.84	0.75	0.92	0.97	0.81	Reference
2756226	1956	Ernst Brandl, Hans Margreiter	Oral Penicillin	7	0.71	0.69	0.62	0.76	0.91	0.67	Reference
2797183	1957	Hazen/ Brown	Nystatin	13	0.90	0.57	0.79	0.95	0.84	0.85	Reference
2816721	1957	R. J. Taylor	Rocket Engine	25	0.71	0.90	0.91	0.74	0.98	0.95	Reference
2817025	1957	Robert Adler	TV remote control	27	0.74	0.94	0.92	0.77	0.99	0.95	Reference
2835548	1958	Robert C. Baumann	Satellite	16	0.79	0.87	0.83	0.83	0.98	0.89	Reference
2866012	1958	Charles P. Ginsburg	Video Tape Recorder	30	0.80	0.93	0.93	0.85	0.99	0.96	Reference
2879439	1959	Charles H. Townes	Maser	24	0.73	0.91	0.90	0.77	0.99	0.94	Reference
2929922	1960	Arthur L. Shawlow	Laser	122	0.82	1.00	0.99	0.87	1.00	1.00	Reference
2937186	1960	Burckhalter/Seiwald	Antibody Labelling Agent	8	0.82	0.22	0.66	0.88	0.28	0.72	Reference
2947611	1960	Francis P. Bundy	Diamond Synthesis	62	0.71	0.37	0.98	0.75	0.70	0.99	Reference
2956114	1960	Charles P. Ginsburg	Wideband Magnetic Tape System	11	0.73	0.77	0.75	0.78	0.94	0.81	Reference
2981877	1961	Robert N. Noyce	Semiconductor Device	152	0.95	1.00	1.00	0.98	1.00	1.00	Reference
3057356	1962	Greatbatch Wilson	Pacemaker	127	0.86	0.91	0.99	0.92	0.99	1.00	Reference
3093346	1963	Maxime A. Faget	First Manned Space Capsule-Mercury	19	0.86	0.85	0.86	0.92	0.97	0.91	Reference
3097366	1963	Paul Winchell	Artificial Heart	23	0.45	0.69	0.89	0.36	0.91	0.93	Reference
3118022	1964	Gerhard M. Sessler	Electret Microphone	39	0.69	0.82	0.95	0.73	0.96	0.97	Reference
3156523	1964	Glenn T. Seaborg	Americium (Element 95)	1	0.84	-	0.17	0.90	0.13	0.13	Reference
3174267	1965	Edward C Bopf, Deere & Co	Cotton Harvester	4	0.43	0.57	0.46	0.32	0.84	0.47	Reference
3220816	1965	Alastair Pilkington	Manufacture of Flat Glass	25	0.77	0.37	0.91	0.83	0.69	0.94	Reference
3287323	1966	Stephanie Kwolek, Paul Morgan	Kevlar	1	0.63	-	0.17	0.69	0.08	0.12	Reference
3478216	1969	George Carruthers	Far-Ultraviolet Camera	3	0.63	0.37	0.38	0.82	0.70	0.39	Reference
3574791	1971	Patsy Sherman	Scotchguard	81	0.54	0.84	0.98	0.79	0.97	0.99	Reference
3663762	1972	Edward Joel Amos Jr	Cellular Telephone	112	0.59	0.93	0.99	0.84	0.99	1.00	Reference
3789832	1974	Raymond V. Damadian	MRI	59	0.44	0.89	0.97	0.81	0.98	0.98	Reference
3858232	1974	William Boyle	Digital Eye	51	0.38	0.97	0.97	0.76	1.00	0.98	Reference
3906166	1975	Martin Cooper	Cellular Telephone	219	0.39	0.93	1.00	0.78	0.99	1.00	Reference
4136359	1979	Stephen Wozniak, Apple	Microcomputer for use with video display	37	0.77	0.69	0.95	0.95	0.87	0.94	Reference
4229761	1980	Valerie Thomas	Illusion Transmitter	3	0.84	-	0.38	0.97	0.05	0.21	Reference
4237224	1980	Boyer/Cohen	Molecular chimeras	301	1.00	1.00	1.00	1.00	1.00	1.00	Reference

Table A.1: Important Patents (cont)

Patent	Year	Inventor	Invention	Citations	Percentile Ranks						Source
					No Adjustment			Remove year FE			
					Quality	Citations	Quality	Citations	Quality	Citations	
					(total)	(0-10)	(0-10)	(total)	(0-10)	(0-10)	
4363877	1982	Howard M. Goodman	Human Growth Hormone	51	1.00	0.85	0.97	1.00	0.95	0.96	Reference
4371752	1983	Gordon Matthews	Digital Voice Mail System	223	0.82	0.98	1.00	0.92	1.00	1.00	Reference
4399216	1983	Richard Axel	Co-transformation	482	0.99	0.98	1.00	1.00	1.00	1.00	Reference
4437122	1984	Walsh/Halpert	bitmap (raster) graphics	178	1.00	0.96	1.00	1.00	0.99	1.00	Reference
4464652	1984	Apple	Lisa Mouse	112	0.85	0.98	0.99	0.92	1.00	0.99	Reference
4468464	1984	Boyer/Cohen	Molecular chimeras	109	1.00	0.91	0.99	1.00	0.97	0.99	Reference
4590598	1986	Gordon Gould	Laser	20	0.76	0.49	0.87	0.62	0.33	0.80	Reference
4634665	1987	Richard Axel	Co-transformation	183	1.00	0.82	1.00	1.00	0.88	1.00	Reference
4683195	1987	Kary B. Mullis	polymerase chain reaction	2884	1.00	1.00	1.00	1.00	1.00	1.00	Reference
4683202	1987	(several)	polymerase chain reaction	3328	0.99	1.00	1.00	0.99	1.00	1.00	Reference
4736866	1988	Leder/Stewart	transgenic (genetically modified) animals	370	1.00	0.99	1.00	1.00	1.00	1.00	Reference
4744360	1988	Patricia Bath	Cataract Laserphaco Probe	81	0.94	0.97	0.98	0.89	0.99	0.98	Reference
4816397	1989	Michael A. Boss	recombinant antibodies	567	0.99	0.97	1.00	0.99	0.99	1.00	Reference
4816567	1989	Shmuel Cabilly	immunoglobulins	1785	0.99	0.98	1.00	0.99	0.99	1.00	Reference
4838644	1989	Ellen Ochoa	Recognizing Method	22	0.96	0.88	0.89	0.94	0.92	0.81	Reference
4889818	1989	(several)	polymerase chain reaction	366	1.00	1.00	1.00	1.00	1.00	1.00	Reference
4965188	1990	(several)	polymerase chain reaction	1176	1.00	1.00	1.00	1.00	1.00	1.00	Reference
5061620	1991	(several)	Method for isolating the human stem cell	252	0.99	1.00	1.00	0.99	1.00	1.00	Reference
5071161	1991	Geoffrey L Mahoon	Airbag	23	0.87	0.94	0.89	0.68	0.96	0.81	Reference
5108388	1992	Stephen L. Troke	Laser Surgery Method	125	0.97	0.94	0.99	0.95	0.95	0.99	Reference
5149636	1992	Richard Axel	Co-transformation	6	1.00	0.49	0.58	1.00	0.22	0.36	Reference
5179017	1993	Richard Axel	Co-transformation	131	1.00	0.95	0.99	1.00	0.96	0.99	Reference
5184830	1993	Saturo Okada, Shin Kojo	Compact Hand-Held Video Game System	201	0.98	0.99	1.00	0.97	1.00	1.00	Reference
5194299	1993	Arthur Fry	Post-It Note	76	0.89	0.80	0.98	0.69	0.78	0.97	Reference
5225539	1993	Gregory P. Winter	Chimeric, humanized antibodies	671	1.00	1.00	1.00	1.00	1.00	1.00	Reference
5272628	1993	Michael Koss	Core Excel Function	94	1.00	0.98	0.99	0.99	0.99	0.98	Reference
5747282	1998	Mark H. Skolnick	isolating BRCA1 gene	15	0.95	0.69	0.82	0.91	0.32	0.67	Reference
5770429	1998	Nils Lonberg	human antibodies from transgenic mice	248	0.83	0.99	1.00	0.46	0.99	1.00	Reference
5837492	1998	(several)	isolating BRCA2 gene	5	0.82	0.37	0.52	0.44	0.08	0.26	Reference
5939598	1999	(several)	Transgenic mice	262	1.00	0.94	1.00	1.00	0.93	1.00	Reference

Table A.1: Important Patents (cont)

Patent	Year	Inventor	Invention	Citations	Percentile Ranks						Source
					No Adjustment			Remove year FE			
					Quality	Citations		Quality	Citations		
					(total)	(0-10)	(0-10)	(total)	(0-10)	(0-10)	
5960411	1999	Peri Hartman, Jeff Bezos	1-click buying	1387	1.00	1.00	1.00	1.00	1.00	1.00	Reference
6230409	2001	Patricia Billings	Geobond	7	0.77	0.69	0.62	0.74	0.36	0.46	Reference
6285999	2001	Larry Page	Google Pagerank	689	0.98	1.00	1.00	0.99	1.00	1.00	Reference
6331415	2001	Shmuel Cabilly	Antibody molecules	243	1.00	0.49	1.00	1.00	0.18	1.00	Reference
6455275	2002	Richard Axel	Co-transformation	7	0.93	0.49	0.62	0.97	0.19	0.52	Reference



Table A.2: Validation: Patent Importance and Forward Citations

Forward Citations	A. Contemporaneous Relation						B. Predictive Relation					
	(0, 1) → (0, 1)		(0, 5) → (0, 5)		(0, 10) → (0, 10)		(0, 1) → 2+		(0, 5) → 6+		(0, 10) → 11+	
Horizon	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)	(11)	(12)
log(Patent Importance)	0.275*** (0.062)	0.139 (0.074)	1.008*** (0.037)	0.789*** (0.053)	1.058*** (0.015)	0.894*** (0.025)	1.029*** (0.059)	0.997*** (0.064)	0.692*** (0.053)	0.769*** (0.075)	0.344*** (0.027)	0.391*** (0.045)
log(1 + Fwd. Citations )							0.615*** (0.017)	0.512*** (0.015)	0.588*** (0.016)	0.550*** (0.015)	0.546*** (0.015)	0.517*** (0.014)
$R^2$	0.092	0.225	0.232	0.367	0.295	0.425	0.354	0.508	0.362	0.497	0.347	0.472
Observations	6,017,673	4,084,292	4,802,836	3,064,631	4,135,358	2,533,724	6,017,673	4,084,292	4,964,003	3,195,838	4,135,358	2,533,724
Grant Year FE	Y		Y		Y		Y		Y		Y	
Technology Class (CPC3)	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y
Assignee × Year FE		Y		Y		Y		Y		Y		Y

Table reports the results of estimating the following specification at the patent level (indexed by  $j$ ):

$$\log(1 + CITES_j) = \alpha + \beta \log q_j^\tau + \gamma \mathbf{Z}_j + \varepsilon_j.$$

In terms of the independent variables, we measure patent importance and citations over the  $\tau$  years since the patent is filed. For the dependent variable, we measure forward citations over the same interval (Panel A) or from year  $\tau + 1$  onwards (Panel B). The vector  $\mathbf{Z}_j$  includes dummies controlling for technology class (defined at the 3-digit CPC level), grant year, and the interaction of assignee and year effects. Including assignee fixed effects reduces the number of observations since many patents have no assignees. We restrict attention to the sample of patents issued after 1947, as this is the period for which citations are recorded consistently by the USPTO. We cluster the standard errors by the patent grant year and report them in parentheses. \*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

Table A.3: Validation: Patent Importance and Market Values

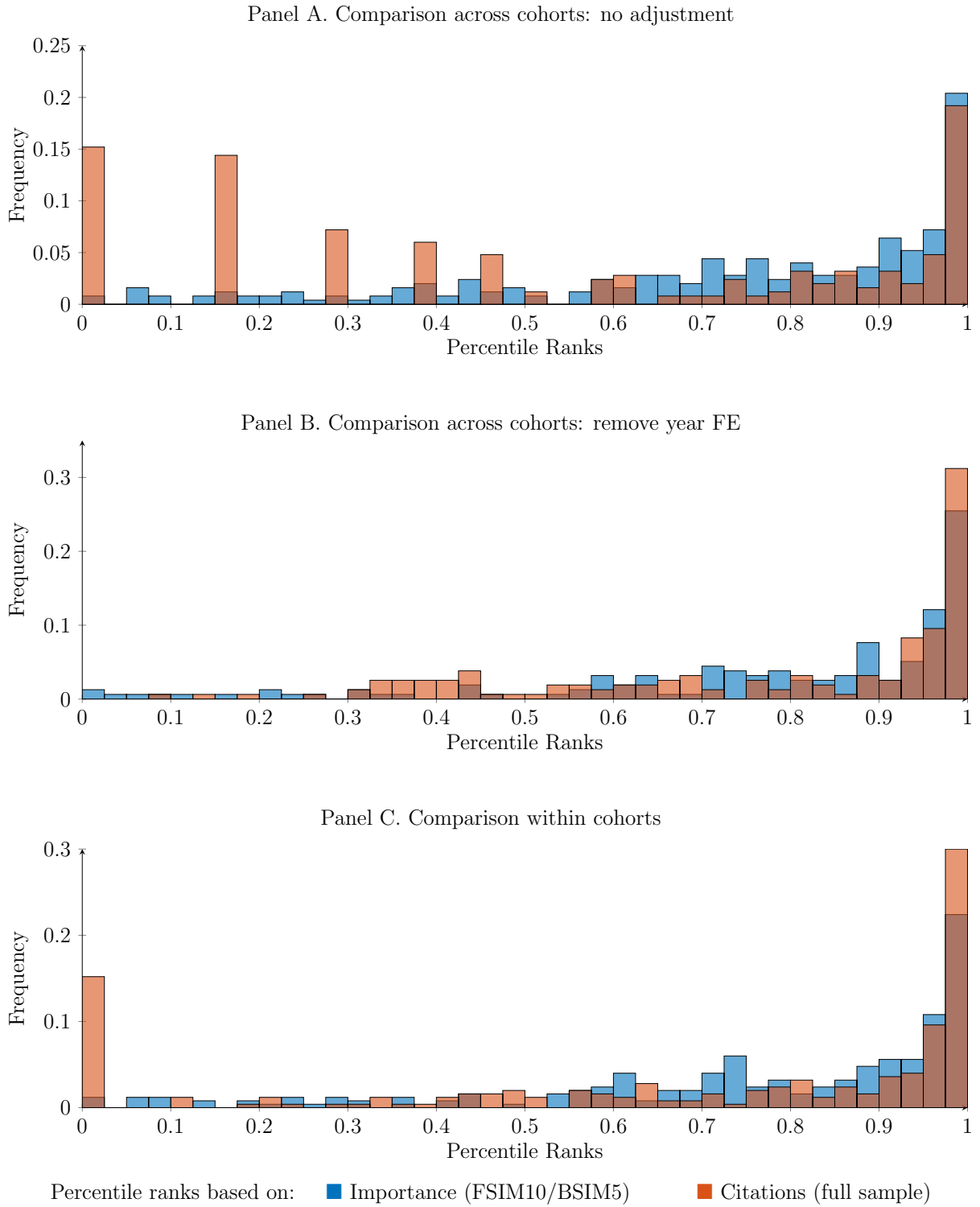
Horizon of Forward Similarity/Citations	(0-1yrs)		(0-5yrs)		(0-10yrs)	
	(1)	(2)	(3)	(4)	(5)	(6)
log(Patent Importance)	0.0019** (0.0009)	0.0020** (0.0010)	0.0027** (0.0012)	0.0024* (0.0012)	0.0039** (0.0015)	0.0029* (0.0016)
log(1 + Forward Citations )		-0.0003 (0.0004)		0.0016** (0.0006)		0.0038*** (0.0010)
<i>Observations</i>	2,097,976	2,097,976	1,737,721	1,737,721	1,424,918	1,424,918
<i>R</i> <sup>2</sup>	0.949	0.949	0.947	0.947	0.939	0.939

Table reports the results of estimating the following specification

$$\log \hat{V}_j = \alpha + \beta \log q_j^r + \gamma \mathbf{Z}_j + \varepsilon_j.$$

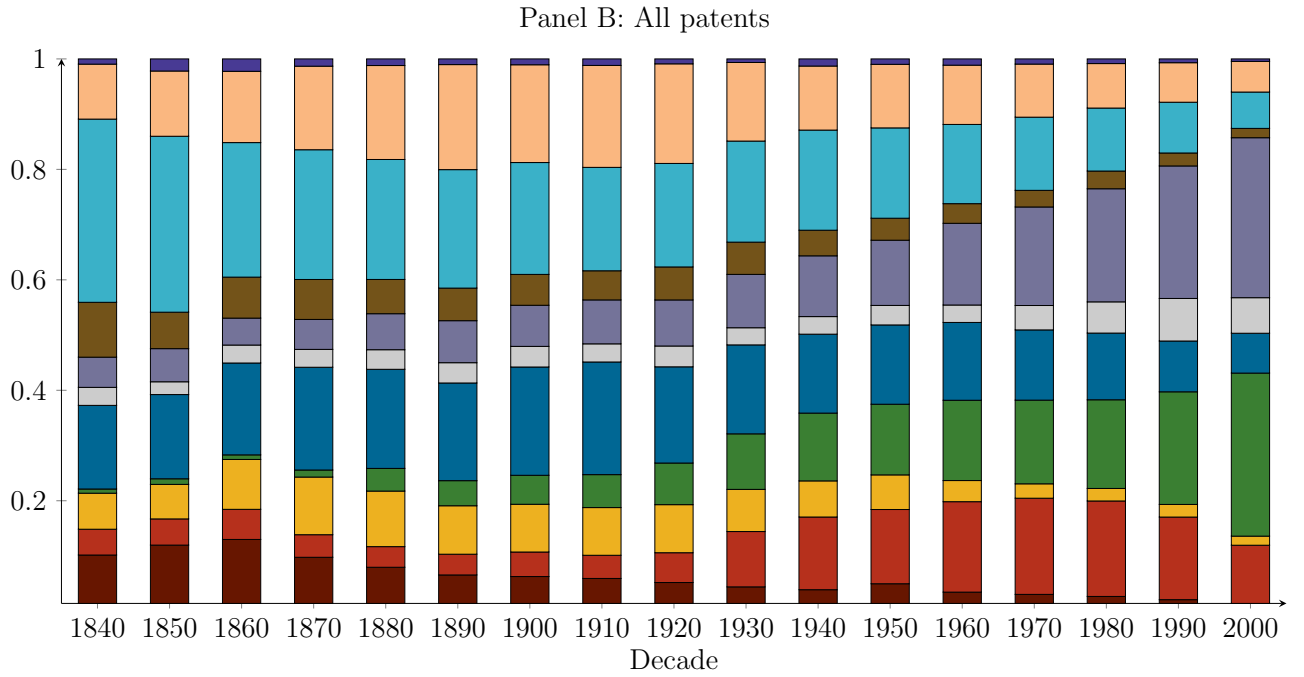
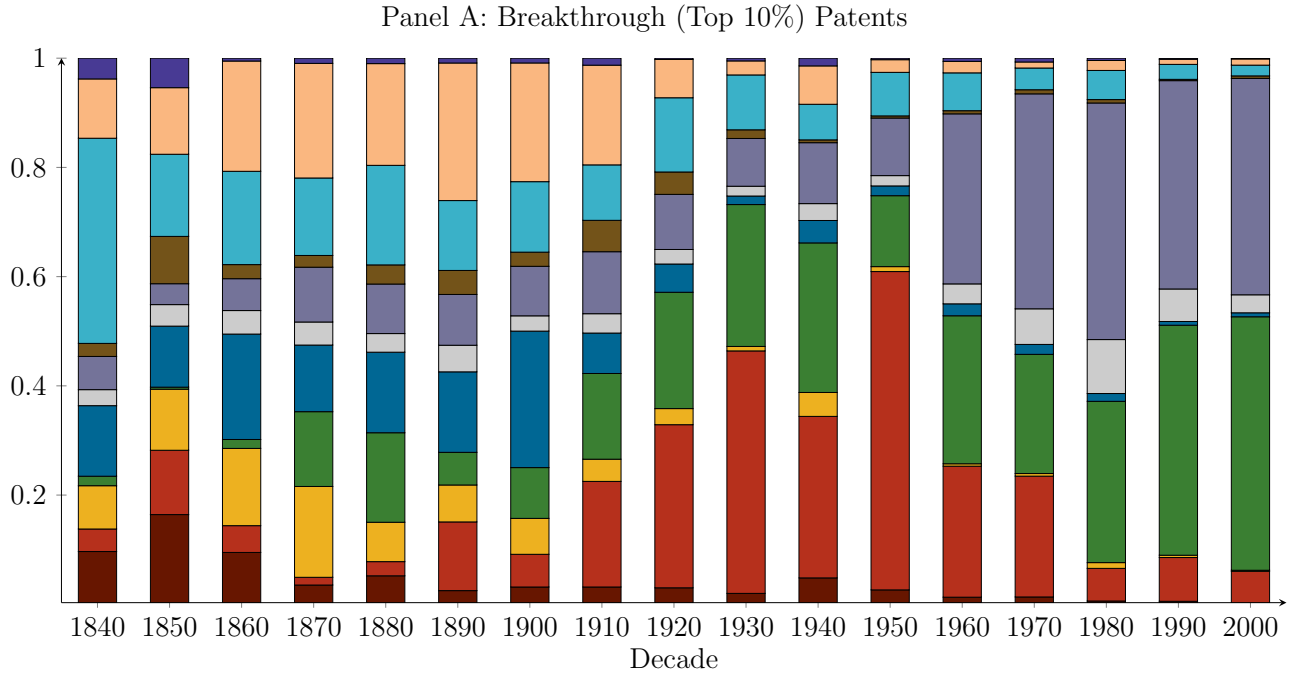
The regression relates the log of the Kogan et al. (2017) estimate of the market value of the patent to our (log) measures of patent importance, which combines the patent's impact and novelty, constructed in equation (10). As controls  $\mathbf{Z}_j$ , we include dummies controlling for technology class (defined at the 3-digit CPC level), the logarithm of the firm's market capitalization and the interaction of firm (CRSP: permco) and grant year effects. In columns (3), (5), and (5) we include as additional controls the number of forward citations (measured over the same horizon as our importance measure). We cluster the standard errors by the patent grant year and report them in parentheses. Independent variables are normalized to unit standard deviation. \*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

**Figure A.1: Significant Patents: Importance vs Forward Citations**



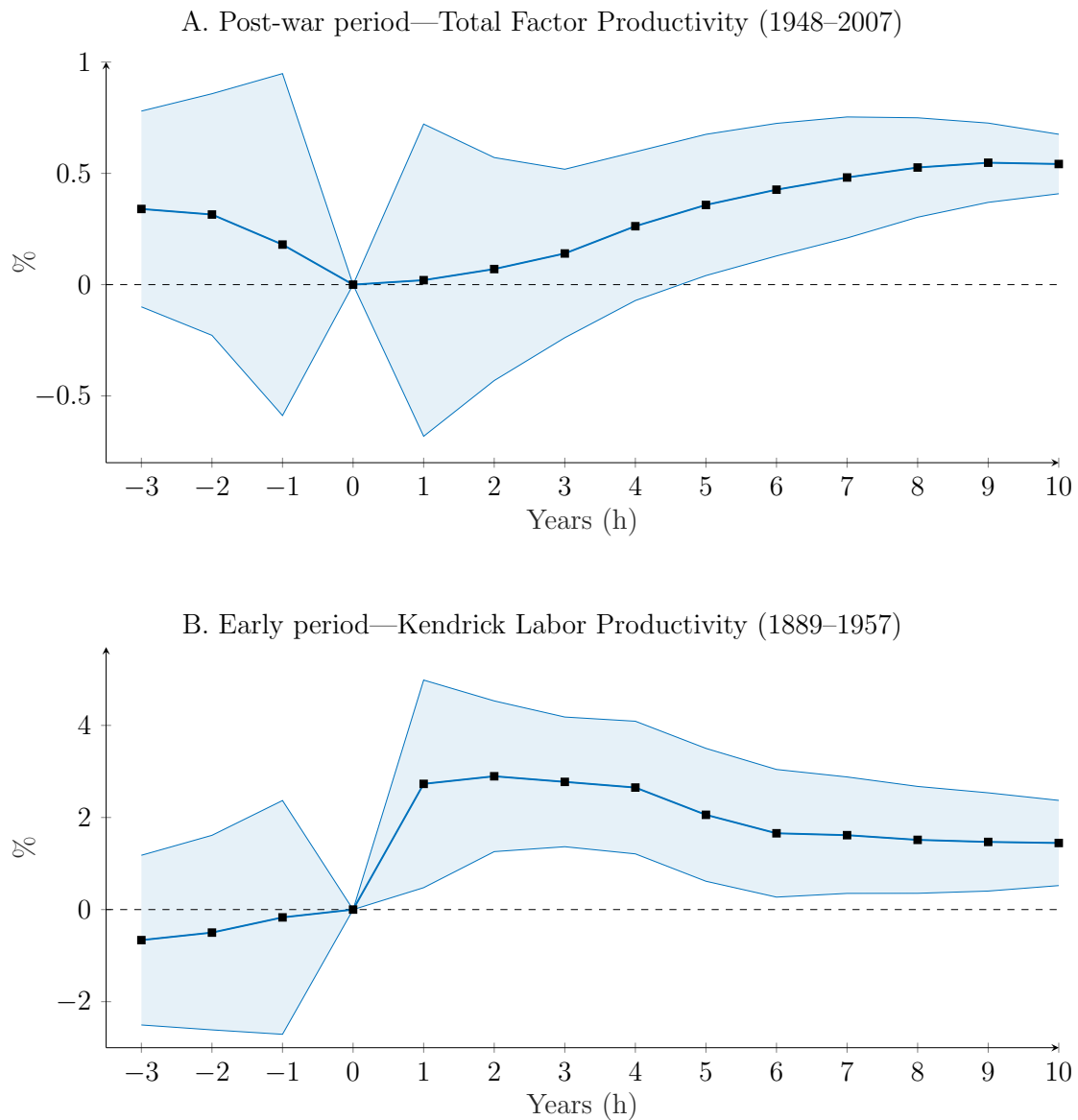
Distribution of patent percentile ranks based on our patent importance indicator (blue) measured over a horizon of 10 years and forward citations (light red) measured over the entire sample. A value of x% indicates that a given patent scores higher than x% of all other patents in the sample (panel A); same after removing year-fixed effects from importance and citations (Panel B); or computing percentile ranks relative to patents that are issued in the same year (panel C). The list of patents, along with their source, appears in Appendix Table A.1

Figure A.2: Breakdown of Innovation by Technology Classes



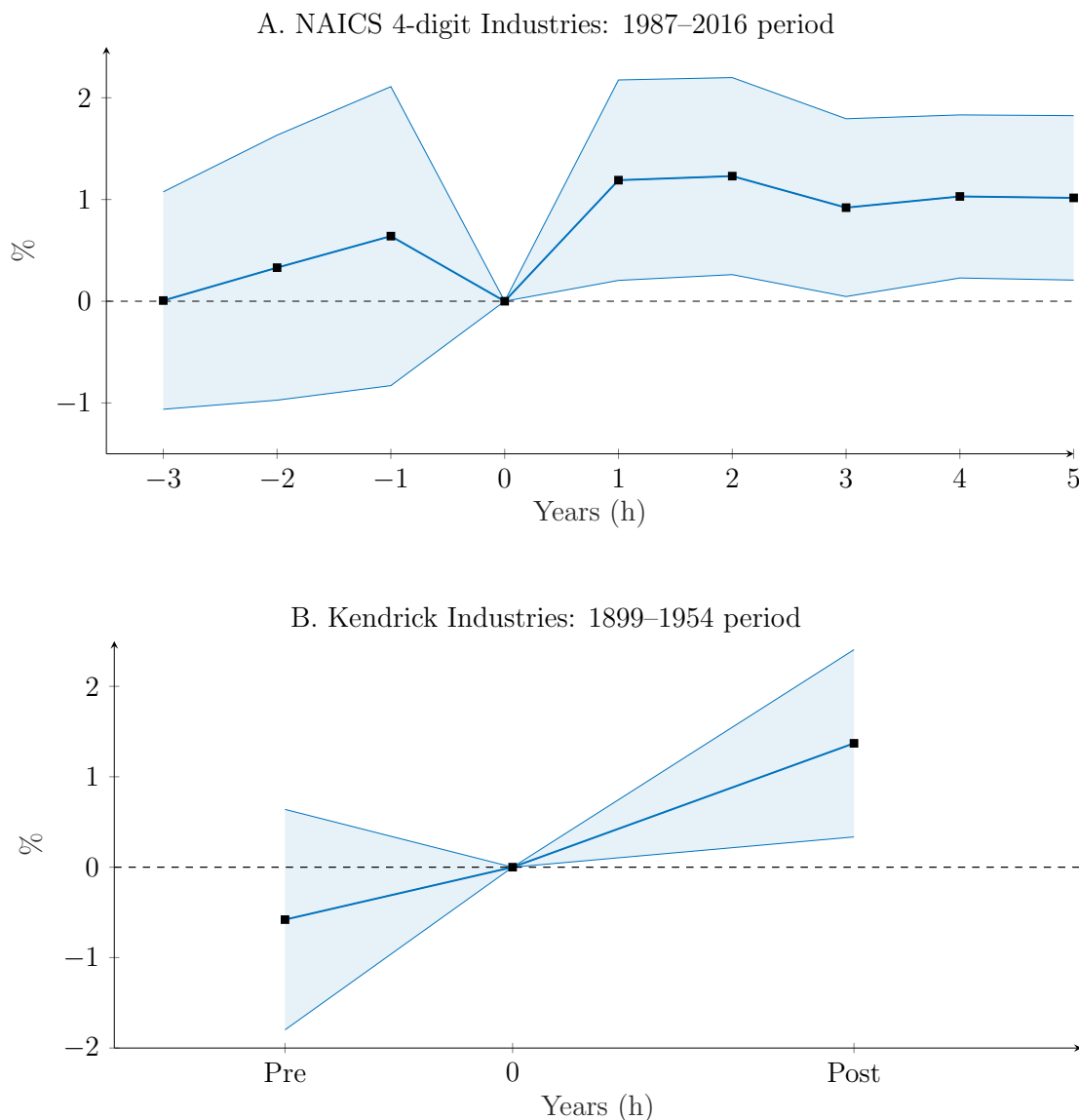
- Agriculture and Food (A0, A2)
- Electricity and Electronics (H0)
- Health and Entertainment (A6)
- Lighting, Heating, Nuclear (F2, G2)
- Transportation (B6)
- Chemistry and Metallurgy (C)
- Engineering, Construction, and Mining (E0, E2, F0, F1)
- Instruments, Information (G, Y1)
- Manufacturing Process (B0, B2, B3, B4, B8, D0, D1, D2)
- Weapons (F4)
- Consumer Goods(A4)

**Figure A.3: Breakthrough Innovation and Aggregate TFP**



Response of measured productivity to a unit standard deviation shock to our technological innovation index (in logs). In Panel A, productivity is measured using total factor productivity from Basu et al. (2006). In Panel B, productivity is measured by output per manhour in manufacturing (Kendrick, Table D-II). We include 90% confidence intervals, computed using Newey-West standard errors (with a maximum number of lags equal to one plus the number of overlapping observations). All specifications control for the lag level of productivity.

**Figure A.4: Breakthrough Innovation and Industry TFP**



Response of industry total factor productivity to a unit standard deviation shock to our technological innovation index. Panel A presents results for 86 manufacturing industries at the NAICS 4-digit level. Productivity data is from the Bureau of Labor Statistics. Kendrick industries are from Table D-V, and productivity is output per manhour. The Kendrick data includes information for the level of labor productivity (output per manhour) for 62 manufacturing industries for the years 1899, 1909, 1919, 1937, 1947, and 1954. For each period ( $t, s$ ), we regress the annualized difference in log labor productivity on the log of the accumulated level of innovation (number of breakthrough patents) in  $t \pm 2$  years—controlling for time and industry dummies, the log number of patents during the same period, and the log level of productivity at  $t$ . Standard errors are clustered by industry. To construct industry innovation indices for NAICS industries, we use the probabilistic mapping from CPC codes to NAICS codes from Goldschlag et al. (2016). To construct innovation indices for the Kendrick industries, which are defined at the SIC code level, we use the concordance between 1997 NAICS and 1987 SIC codes from the Census Bureau. If NAICS industries map into multiple SIC codes, we assign an equal fraction to each.

**Figure A.5: Breakthrough patents and Industry TFP—comparison to Citations**

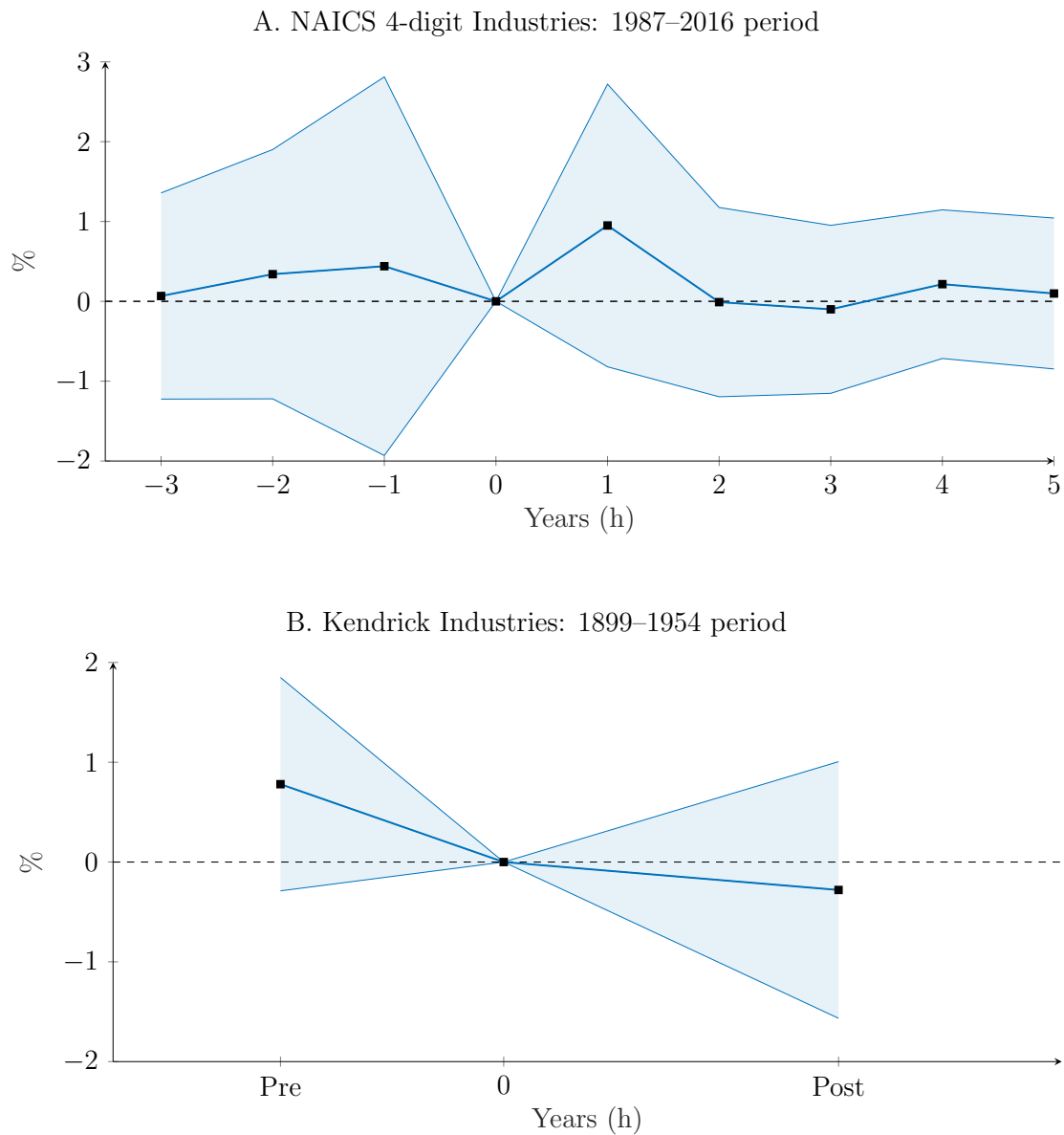


Figure performs the same exercise as Figure A.4, except that we now construct the industry innovation indices based on citation counts.