

# Technology and Real Options: Evidence from Patent Text

PRELIMINARY AND INCOMPLETE

Davidson Heath\*

November 9, 2014

## **Abstract**

I map U.S. firms' technological positions using textual analysis of their patents. Firms' technological positions positively predict their product market positions. Greater technological differentiation is associated with higher market value, profitability and residual productivity. Consistent with a model of real options on heterogeneous assets, firms that are more technologically differentiated have lower stock returns and this effect is concentrated in growth firms.

---

\*USC Marshall School of Business. Email: [dheath@usc.edu](mailto:dheath@usc.edu). Many thanks to my committee members Wayne Ferson (chair), Scott Joslin, Gordon Phillips, Gerard Hoberg, and Kenneth Ahern for their advice and encouragement.

# 1 Introduction

Technological innovation is a large, persistent component of productivity growth at both the aggregate and firm level (Solow (1962); Romer (1990)). In this paper I construct new measures of innovation using the text of all 6.2 million utility patents granted by the U.S. Patent and Trademark Office (USPTO) from 1926 to 2010. These new measures open the “black box” of simply counting patents and present a more detailed picture of the nature of innovation at the patent and firm level. I use these measures to estimate the connections between innovation and firms’ productivity and stock returns.

While it is well accepted that firms differ in terms of their product market positions – for example, the industries in which they operate – technology is still generally treated as a homogeneous good that firms amass. Instead, I map patents and firm-years to vectors in “word space”. At the firm-year level, I measure firms’ technological differentiation relative to other firms. I find that technological differentiation is strongly positively associated with average  $Q$  and unexplained productivity residuals (TFP). These relations persist after controlling for other known correlates of firm productivity and value, including measures of market power and product market differentiation and after sweeping out firm fixed effects. Relative to standard measures of technology inputs and outputs, the text-based measures contribute substantial explanatory power.

I compare my patent text-based measures in detail with those of Hoberg and Phillips (2010), who locate firms in product market space based on key words in the business description section of their 10-K filings. There is substantial overlap between Hoberg and Phillips (2010)’s product market positions and my technological positions, but there is also considerable discrepancy which appears to be meaningfully due to differences between firms’ patenting activities and their operating activities. I find that technology similarity positively

predicts future product market similarity, but the reverse is not the case.

Technological differentiation is one way that firms can achieve product market differentiation, by inventing new products. Differentiation in product markets produces less volatile cash flows and reduces firms' risk and expected returns. This effect is magnified when the product is not yet implemented, because real options are levered relative to assets in place (Berk, Green and Naik (2004)). I build a simple model in which innovation yields real options to introduce new products. In the model, while firms' market to book ratio summarizes the risk of their product market positions, it does not summarize the risk of their options on future products. Consistent with the model, I find that technological differentiation predicts lower returns while measures of product market differentiation do not, and the return predictability is concentrated in growth firms.

## 1.1 Prior Literature

In Pastor and Veronesi (2003) and Garleanu, Panageas and Yu (2012), a new technological era arrives and is gradually adopted in the economy. Here, I suppose that firms are positioned in a multidimensional "technology space" that is a precursor to their position in product market space cf. Hotelling or Hoberg-Phillips.

This paper joins other recent studies that apply textual analysis to measure aspects of firm behavior beyond the standard accounting variables. The closest related paper is Hoberg and Phillips (2010) who locate firms in product market space based on key words in the business description section of their 10-K filings. Aside from being based on different sources, my measures differ from theirs in two main ways. First, the quantity of text is much larger. The typical business description section of a firm's annual 10-K is less than a page, and sometimes just a few sentences. The body text of a U.S. patent is on average 22 pages

long, and firms often file multiple patents per year. Second, [Hoberg and Phillips \(2010\)](#)'s TNIC measures are available for publicly traded firms from 1997 onward. By contrast, I collect the text of all U.S. patents granted to individuals, public firms, and private firms since 1961. Section 3.3 compares my measure to theirs in detail.

[Packalen and Bhattacharya \(2012\)](#) analyze U.S. patent text: their focus is on the changing nature of technology over time, and they do not link their text based measures to the CRSP/Compustat data. [Alexopoulos \(2011\)](#) finds that innovative activity as measured by the publication of computer- and telecommunications-related technical manuals is followed by growth in GDP, TFP, investment and hours worked at the aggregate level.

[Hirshleifer, Hsu and Li \(2012\)](#); [Hirshleifer, Hsu and Li \(2013\)](#); [Cohen, Diether and Malloy \(2013\)](#) find predictable patterns in stock returns from firms' patenting activity. They find that firms with diverse patents, high-quality patents and high research productivity (patents per R&D dollar) have higher subsequent returns that are not captured by standard return benchmarks. The main interpretation in all these papers is that the market is slow to recognize innovative firms' superiority, so the documented return spreads represent mispricing. By contrast, I find that firms with more technological differentiation have *lower* subsequent returns, and I offer a simple rational explanation based on the leverage effect of growth options on the risk premiums of the underlying assets.

## 2 Data and Methodology

### 2.1 Patent Text

I use the full text of all U.S. utility patents granted from 1926 to 2010 from the U.S. Patent and Trademark Office, which is cross hosted at the Google Patents project<sup>1</sup>. Figure 1 plots the number of patents granted each year that appear in the sample, compared to the official grant numbers by year from the USPTO which are available starting in 1963. The sample covers well over 99% of all U.S. utility patents every year from 1963 to 2009, with a coverage ratio of 85% in 2010. The sample comprises 6,237,597 patents and 240GB of text, which corresponds to roughly 134 million pages.

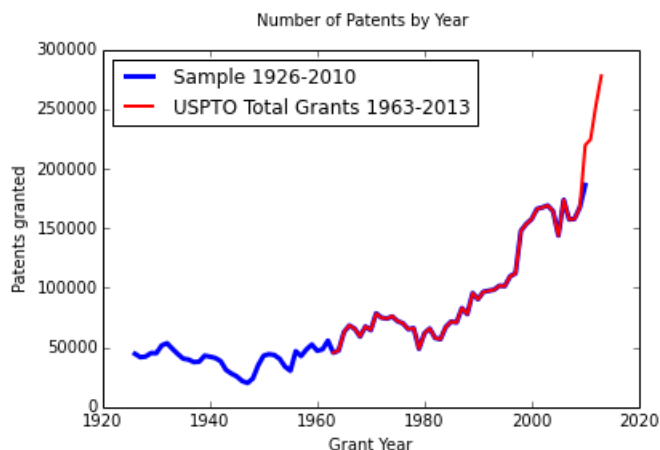


Figure 1: Number of patents each year in the sample from 1926-2010 and number of patents granted by the USPTO from 1963-2013.

For each patent I take the full body text and extract all words (1-grams) and two and three word phrases (2-grams and 3-grams). I select words and phrases made up of only uncapitalized English nouns that appear in no fewer than 1000 and no more than 1% of all

---

<sup>1</sup><http://www.google.com/googlebooks/uspto-patents-grants-text.html>

patents in the sample (that is, no more than 62,238 patents). This yields a list of 64,369 nouns and noun phrases, or 'ngrams'. The details of how I filter and select the ngrams are in Appendix A. My results are similar when I experiment with different rules for selecting the ngram list.

Table 1 lists the most common ngrams by their vintage (the application year of the first patent that mentions the ngram) for the years 1950, 1960, and so on up to 2000, which illustrates the changes in technology over the period. The vintages support the conclusion of Griliches (1990) that patents are applied for early in the inventive process. “GPS system” first appears in a patent in 1980; the first functioning GPS satellite was not launched until 1989. “Notebook computer” is first mentioned in 1990 when such portable computers were quite unusual. “Picture experts” and “picture experts group”, which first appear in a patent in 1990, refer to the MPEG digital video format, which was first discussed in 1988 and the first MPEG format was officially released in 1993.

Table 1: The table shows the most common ngrams (one, two or three word noun phrases) by their vintage (year of first mention) for vintages 1950, 1960, 1970, 1980, 1990 and 2000.

1950	1960	1970
control circuitry	system memory	acid sequences
clock cycle	glass transition temperature	methodologies
substrate material	status information	software components
epoxy resins	immune response	bus interface
computer networks	control module	programming languages
command signal	calf serum	interface card
remote computer	fibroblast	plasma display
enantiomers	error message	browsers
breast cancer	disclosure bulletin	management information
polymer matrix	logic device	cpu controls

1980	1990	2000
program product	email address	communication media computer
cytokines	notebook computer	wireless media combinations
fusion protein	remote memory storage	computer communication media
protein expression	notebook computers	laptop devices multiprocessor
flash chromatography	email messages	systems methods features
hybridization conditions	picture experts	methods components materials
gps system	picture experts group	access protocol soap
angioplasty	email addresses	amplification markers
necrosis factor	sound card	server data stores
cdna clones	multiple servers	valdecoxib

For each patent, I evaluate if each ngram is present (1) or absent (0), yielding a Boolean vector for each patent. I summarize patents in this way because patents vary greatly in their length and structure – mechanical patents often consist mostly of diagrams, while biotech patents often run to hundreds of pages. Thus, above the simple presence or absence of a given ngram, the number of times the ngram is mentioned in a patent is uninformative.

Some patent documents are almost entirely diagrams of the invention. I exclude patents

from the analysis if they contain fewer than 10 ngrams, as these are likely more noise than signal. About 46,000 patents fall into this category, almost all of them from the early years of the sample. After these filters, the word space is populated by 6.2 million vectors representing the vast majority of U.S. utility patents granted from 1926 to 2010.

Appendix B describes the construction of quality and vintage scores at the level of individual patents. Briefly, I find that patents' quality and vintage positively predict whether the patent holders pay fees to renew them over and above standard measures of quality like citations. This suggests that my ngram selection and vectorization are capturing meaningful information about the patents.

### 3 Firm-year Technology Positions

#### 3.1 Firm Data

Annual firm data is from Compustat. I drop utility, financial and government firms (firms whose SIC codes begin with 49 or 6 or 9) from the sample. I drop firm-years with missing or negative assets or sales and winsorize all ratios at the 1% level in both tails.

I map patents to firms using the database of [Kogan et al. \(2014\)](#), which builds on and extends the NBER patent database ([Hall, Jaffe and Trajtenberg \(2001\)](#)). For the results in this section I index patents by their application year, as is standard in the literature. The sample of firm-years starts in 1961. Because I have data on patents that were granted up to 2010, I end the sample in 2007, allowing a three year lag from application to grant. I compute two standard innovation measures for each firm-year: the count of patents and the count of subsequent citations to those patents, adjusted by category and year as in [Hall, Jaffe and Trajtenberg \(2001\)](#).



The main departure from the full Compustat sample is that my text-based measures require a firm has at least one patent to its name. Hence, firms appear in the panel in the year that they first applied for a patent that was subsequently granted. Table 2 presents summary statistics for the sample compared to the full Compustat cross section. Patenting firms represent 25-45% of all Compustat firms in any year, and the overall distribution of firm characteristics is similar between the subsample and the full set. Patenting firms are larger, higher average  $Q$ , and older. They have similar investment as a fraction of book assets, and are more  $R\&D$  intensive relative to the full set of all Compustat firms.

Table 2: Summary statistics of the Compustat firm-years sample.

	<u>Patenting Firms 1961-2007</u>					<u>All Compustat Firms 1961-2007</u>				
	Mean	Median	St. Dev.	$p_{10}$	$p_{90}$	Mean	Median	St. Dev.	$p_{10}$	$p_{90}$
log(AT)	5.67	5.45	1.90	3.33	8.27	4.94	4.66	1.74	2.88	7.40
Average $Q$	1.93	1.38	1.64	0.83	3.60	1.74	1.27	1.44	0.80	3.10
Firm Age	12.8	10	10.6	1	29	10.3	8	9.5	1	24
CAPX/AT	0.065	0.053	0.050	0.018	0.125	0.072	0.051	0.071	0.013	0.155
R&D/AT	0.059	0.026	0.092	0.000	0.152	0.034	0.000	0.078	0.000	0.108
Firms per Year	858.2	849	290.9	429	1295	2922.8	3117	1223.6	956	4687

## 3.2 Firm-year Technology Vectors

For each firm-year in the data, I sum the vectors of the patents that are mapped to the firm in that year. Thus, the entries of the firm-year vectors are integer counts of how many firm patents contained each ngram. This is in contrast to the approach for individual patents where all that matters is if an ngram is present (1) or absent (0). At the firm-year level, each patent represents a separate invention, so the number of the firm's patents that mention a given ngram is informative about the firm's research program.

Small firms may have several patents one year and zero the next. I fill the resulting gaps in the panel in two ways. The first way is to construct word stock vectors: these vectors are a rolling stock of patent counts from the focal year plus the stock from the prior year, 'depreciating' the prior year by some amount  $\delta \geq 0$ . The second way is to fill the firm-year vectors forward: when a firm had no patents in some year, but did have patents in a prior year, I copy the most recent prior year's vector. I present the results using word stock vectors with depreciation rate  $\delta = 0.2$ , which is used to construct R&D and patent stocks in previous literature. All results in the paper are robust to using other depreciation rates such as  $\delta = 0.1$  or  $\delta = 0$  or filling forward instead.

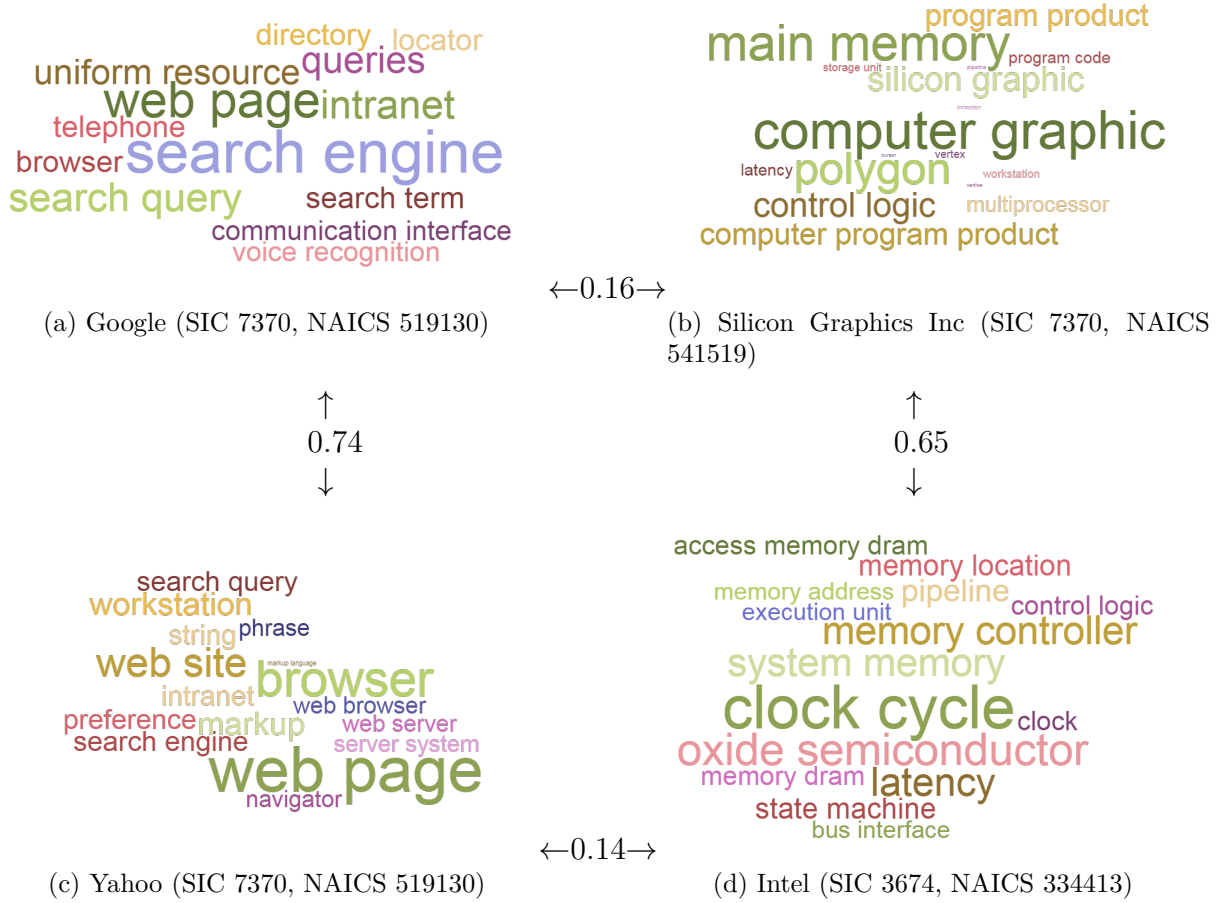


Figure 2: Word clouds for Google, Yahoo, SGI and Intel in 2000. The clouds display the twenty ngrams that appear in the largest number of patents that each firm applied for in 2000. The size of the ngram increases with the fraction of the firm’s patents that mention the ngram. Cosine similarity scores between firm-years are shown with arrows.

Figure 2 displays “word clouds” based on the vectors of four firms for the year 2000: Google, Yahoo, Silicon Graphics Inc and Intel. The word clouds display the 20 most common ngrams for those firms and the size of the typeface increases with the fraction of the firm’s patents in which the ngram appeared. The word clouds appear to be reflective of the firms’ technological areas (note that Google was a small private firm as of the year 2000). There is

technological similarity both within and across conventional industries: SGI and Intel have high similarity despite being in different 1-digit SIC and NAICS industries.

### 3.3 Comparison with Hoberg-Phillips TNIC

This section compares my patent-based similarity measures with the [Hoberg and Phillips \(2010\)](#) TNIC classifications, which map firms’ product market position based on their 10-K product descriptions. The conclusion is that 10-K text captures where firms *are*, while patent text captures where firms are *heading*. To compare the overlap between the two data measures, in this section I drop years prior to 1997 from my data and years after 2007 from the Hoberg-Phillips data. I retain only gvkey-years that appear in both data sets, so firm  $i$  in year  $t$  is only present if firm  $i$  was public i.e. filed a 10-K in year  $t$  and if it had at least one patent in year  $t$  or any year prior.

As the raw similarity scores for the Hoberg-Phillips product market network are not available, I change my similarity scores into a neighbor network to facilitate comparison in this section. For each firm-firm pair in each year, I assign them to be technology (“TECH”) neighbors if their cosine similarity based on patent text is higher than 0.17. This cutoff rule produces a network that, distributionally, is a very close match to that of the Hoberg-Phillips product market (“PROD”) network. Figure 3 plots the histograms of the number of TECH neighbors and PROD neighbors that each firm-year has in the data.

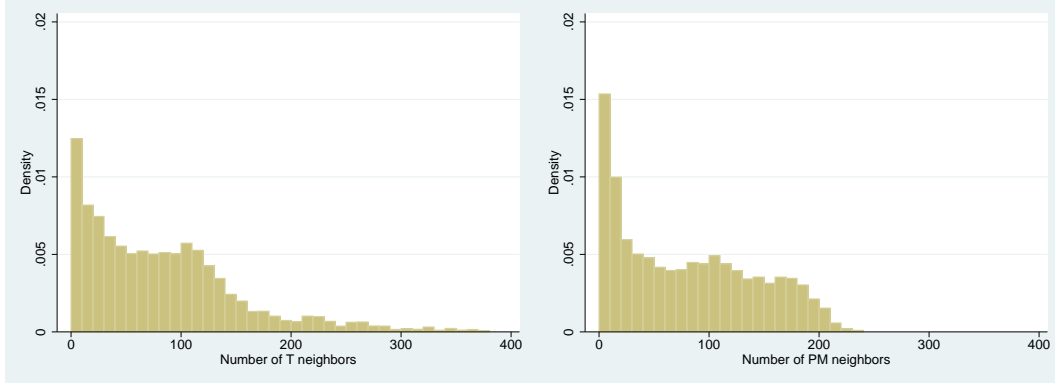


Figure 3: Histograms of the degree (number of network neighbors) of all firm-years for the technology (TECH) and Hoberg-Phillips (PROD) networks. The mean and median degree of the TECH network are 81.8 (69), those of the PROD network are 79.3 (71).

Averaged across the years, any two firms have a 3.7% (3.5%) unconditional chance of being TECH neighbors (PROD neighbours). Conditional on being TECH neighbors, any two firms have a 31.8% chance of also being PROD neighbors; the chance of PROD neighbors also being TECH neighbors is 32.7%. Thus, there is overlap between the two networks far in excess of what we would expect by chance, but firms' TECH and PROD neighbor lists also differ significantly.

Table 3: Logit regressions of product market and technology “neighborhoodness”. The variables are dummies for whether a given firm-firm pair were Hoberg-Phillips neighbors in year  $t$  ( $PRODNEIGHBOR_t$ ) and/or technology neighbors in year  $t$  ( $TECHNEIGHBOR_t$ ). The coefficients report the estimated marginal effects. The standard errors are clustered by year.

	(1) $PRODNEIGHBOR_{t+1}$	(2) $TECHNEIGHBOR_{t+1}$
$TECHNEIGHBOR_t$	0.024*** (0.0075)	0.66*** (0.012)
$PRODNEIGHBOR_t$	0.32*** (0.021)	-0.025*** (0.0018)
Observations	3,461,540	3,461,540
Fixed Effects	Year	Year
Pseudo R-squared	0.14	0.38
Robust standard errors in parentheses		
*** p<0.01, ** p<0.05, * p<0.1		

Table 3 presents the results of logit regressions using dummy variables for whether a given firm-firm pair were product market neighbors and/or technology neighbors in year  $t$ . Column 1 shows that firms that are technology neighbors this year are more likely to become product market neighbors next year. By contrast, firms that are product market neighbors this year are *less* likely to be technology neighbors next year, as shown in Column 2. That is, technology similarity positively predicts future product market similarity, but the reverse is not true.

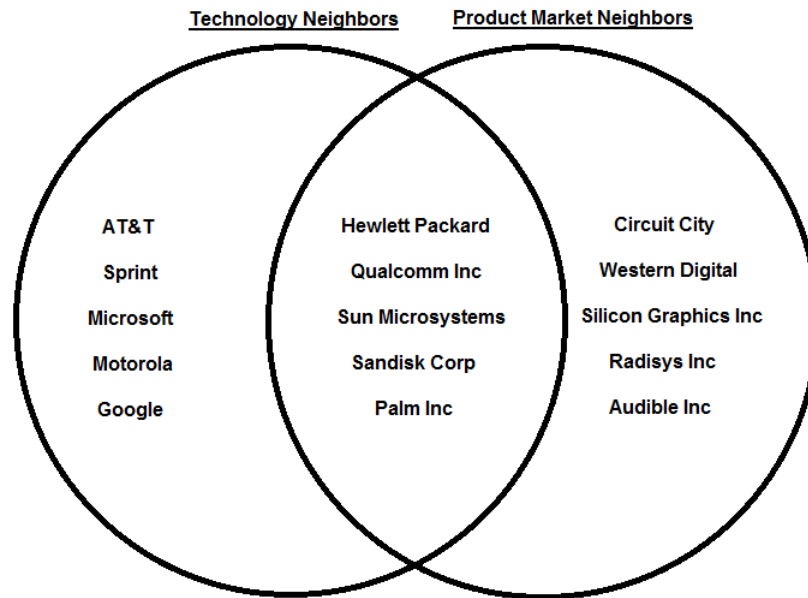


Figure 4: Technology and product market neighbors of Apple Computer in 2006. The technology neighbors are based on similarities in firms' patent text, while the product market neighbors are based on similarities in firms' 10-K product descriptions (Hoberg and Phillips (2010)).

As an illustration, Figure 4 displays the largest five firms that were TECH or PROD neighbors, or both, of Apple Computer in 2006. In 2006, Apple Computer produced personal computers and personal music players (iPods). The long rumored iPhone was announced on January 9, 2007. Several of Apple's patents in 2005 and 2006 dealt with cellular phone technology, with the result that AT&T, Sprint, and Motorola are all TECH-neighbors but not PROD-neighbors. For the same reason, Google was also a TECH-neighbor of Apple in 2006, although neither firm was officially involved in the mobile phone market at the time. In November 2007, Google announced the Android mobile operating system and in 2010 it launched its own line of Nexus smartphones and tablets. Google and Apple first became Hoberg-Phillips (PROD) neighbors in 2010.

### 3.4 Technological Differentiation

Using the firm-year technology vectors, I compute cosine similarities between all firms (public and private) in each year,

$$cos_{ijt} = \left( \frac{firmvec_{it}}{\|firmvec_{it}\|} \right) \cdot \left( \frac{firmvec_{jt}}{\|firmvec_{jt}\|} \right)$$

and for each firm-year define its technological differentiation as:

$$techdiff_{it} = 1 - quantile_{90}(cos_{ijt})$$

, that is, one minus the similarity of firm  $i$  with the firm that is closer than 90% of other firms in the sample that year. The higher the 90% similarity, the closer firm  $i$ 's patent text is to that of other firms in the same area, and the less technologically differentiated it is. The 90% similarities for the firms in Figure 2 are GOOG 0.096, YHOO 0.163, SGI 0.190, INTL 0.245. Thus, Google had the highest *techdiff* and Intel the lowest. Similar results obtain throughout the paper using the mean similarity score for each firm, or a different quantile such as the median. Relative to the median, the 90% similarity score is a more local measure of differentiation.

### 3.5 Intangible Capital

The method of Griliches (1981) is a common means of assessing the value of intangible capital such as R&D or patents (see, e.g., Nicholas (2008)). In the model, firm values take the form

$$V_{it} = \pi_t K_{it} + \eta_t G_{it}$$



where  $\pi$  and  $\eta$  are the shadow prices of tangible and intangible capital respectively. The estimating equation

$$\ln(Q_{it}) = \phi_t G_{it}/K_{it} + \nu_t + \lambda_i + \epsilon_{it}$$

assesses the importance of intangible capital. If a measure of intangible capital  $G_{it}$  is reflected in the firm’s market value but not in its book value then  $\phi_t = \eta_t/\pi_t > 0$ .

Table 4 shows the results of regressing  $\ln(Q)$  by firm-year on *techdiff* plus other common variables. Column 1 shows that both *techdiff* and the log of total patents weighted by adjusted citations are positively related to firm  $Q$  measured within the firm over time (i.e. with firm fixed effects). The size of the coefficient is quite large – a one standard deviation increase in *techdiff* is associated with a 7.8% to 17% higher  $Q$  that year relative to the firm’s sample average.

One concern is if my measures are capturing market power – product market differentiation – as well as new technology. Column 2 adds the Herfindahl of the firm’s SIC industry including the firm itself. The third column adds the Hoberg and Phillips’ text-based product market similarity. The Hoberg-Phillips measures are more granular than SICs and have higher explanatory power for peer firms’ valuation (Hoberg and Phillips (2012)). The coefficient of  $Q$  on *techdiff* remains strongly positive.

Table 4: Measuring intangible assets. The table regresses log of Tobin’s average  $Q$  by firm-year on measures of intangible assets. All independent variables have been standardized to have zero mean and unit variance. The standard errors are clustered by firm and year.

	(1) $\log(Q_t)$	(2) $\log(Q_t)$	(3) $\log(Q_t)$
$techdiff_t$	0.080*** (0.018)	0.078*** (0.018)	0.17*** (0.040)
$\log(cites_t)$	0.014** (0.0059)	0.014** (0.0059)	0.017 (0.011)
$HHI_{SIC,t}$		0.011 (0.0093)	-0.0064 (0.020)
$tnic3tsimm_t$			0.090*** (0.034)
$RDstock/AT_t$	-0.0073 (0.0068)	-0.0072 (0.0068)	-0.059** (0.030)
$\log(AT_t)$	-0.23*** (0.023)	-0.23*** (0.023)	-0.35*** (0.049)
Observations	38,762	38,762	12,946
R-squared	0.041	0.041	0.078
No of firms	4,038	4,038	2,262
Fixed Effects	Firm + Year	Firm + Year	Firm + Year
Robust standard errors in parentheses			
*** p<0.01, ** p<0.05, * p<0.1			

### 3.6 TFP

Imrohoroglu and Tuzel (2014) impute total factor productivity (TFP) by firm-year for Compustat firms using semiparametric methods<sup>2</sup>. They document that high TFP firms tend to be relatively large and relatively high book-to-market (growth) firms. High TFP firms are also younger, less levered, and both their capital and labor grow more rapidly than low TFP firms.

Table 5 regresses firm-year imputed TFP levels on  $techdiff$ , as well as other firm vari-

<sup>2</sup>Many thanks to Profs. Tuzel and Imrohoroglu for sharing their firm-year TFP measurements.

ables that are known to covary with TFP. We see that TFP is strongly positively correlated with technological differentiation. As with firm  $Q$ , the size of the coefficient is again large: a one standard deviation increase in *techdiff* is associated with a 9.4% to 9.8% increase in TFP for that firm-year relative to the firm's sample average.

Column 2 adds the Herfindahl of the firm's SIC industry including the firm itself. The third column adds the Hoberg and Phillips' text-based product market similarity. The measures of product market differentiation are not significantly associated with firm  $TFP$ , and the coefficient on *techdiff* is unchanged.

Table 5: The table regresses TFP imputed by firm-year on the patent text measure *techdiff* plus other control variables that are correlated with firm level productivity. All independent variables have been standardized to have zero mean and unit variance. The standard errors are clustered by firm and year.

	(1) <i>TFP<sub>t</sub></i>	(2) <i>TFP<sub>t</sub></i>	(3) <i>TFP<sub>t</sub></i>
<i>techdiff<sub>t</sub></i>	0.097*** (0.011)	0.098*** (0.011)	0.094*** (0.025)
<i>log(cites<sub>t</sub>)</i>	0.0011 (0.0032)	0.00094 (0.0032)	0.0072 (0.0052)
<i>HHI<sub>SIC,t</sub></i>		-0.012* (0.0068)	-0.0085 (0.019)
<i>tnic3tsimm<sub>t</sub></i>			0.010 (0.025)
<i>RDstock/AT<sub>t</sub></i>	-1.08*** (0.12)	-1.09*** (0.12)	-1.48*** (0.24)
<i>log(AT<sub>t</sub>)</i>	0.27*** (0.020)	0.27*** (0.020)	0.44*** (0.046)
<i>Mkt/Book<sub>t</sub></i>	0.17*** (0.012)	0.17*** (0.012)	0.16*** (0.020)
Observations	30,603	30,603	9,041
R-squared	0.182	0.182	0.180
No of firms	2,974	2,974	1,561
Fixed Effects	Firm + Year	Firm + Year	Firm + Year

Robust standard errors in parentheses

\*\*\* p<0.01, \*\* p<0.05, \* p<0.1

Imputed TFP is a model residual – it is the component of firm revenues or profits that is unexplained by the inputs plus some model of the firm’s production function. Thus, measured TFP depends on the correctness of the model used, in ways that are hard to quantify. If the model used to impute TFP is “more wrong” for innovative firms or industries, then this could drive a spurious correlation between TFP and . In some sense the results above say only that *techdiff* is correlated with “model error”.

Table 6: The table regresses firms' return on assets and investment to assets ratios on the patent text measure *techdiff* plus other control variables that are correlated with firm level productivity. All independent variables have been standardized to have zero mean and unit variance. The standard errors are clustered by firm and year.

	(1) <i>ROA<sub>t</sub></i>	(2) <i>ROA<sub>t</sub></i>	(3) <i>CAPX/AT<sub>t</sub></i>	(4) <i>CAPX/AT<sub>t</sub></i>
<i>techdiff<sub>t</sub></i>	0.030*** (0.0042)	0.038*** (0.0079)	0.0042*** (0.00094)	0.0053*** (0.0015)
<i>log(cites<sub>t</sub>)</i>	-0.00079 (0.0016)	-0.0010 (0.0031)	0.00027 (0.00044)	-0.00044 (0.00067)
<i>HHI<sub>SIC,t</sub></i>	-0.0032 (0.0021)	-0.011* (0.0055)	0.00094 (0.00089)	0.000045 (0.0014)
<i>tnic3tsimm<sub>t</sub></i>		-0.0072 (0.0058)		-0.0023* (0.0014)
<i>RD/sales<sub>t</sub></i>	-0.061*** (0.0052)	-0.073*** (0.010)	0.00035 (0.00053)	-0.00034 (0.00077)
<i>log(AT<sub>t</sub>)</i>	0.095*** (0.012)	0.23*** (0.021)	0.00065 (0.0018)	-0.0011 (0.0028)
<i>Mkt/Book<sub>t</sub></i>	0.016*** (0.0037)	0.023*** (0.0062)	0.0030*** (0.00056)	0.0018*** (0.00042)
Observations	39,957	13,069	39,542	12,994
R-squared	0.134	0.177	0.006	0.005
No of firms	4,080	2,289	4,057	2,278
Fixed Effects	Firm + Year	Firm + Year	Firm + Year	Firm + Year

Robust standard errors in parentheses

\*\*\* p<0.01, \*\* p<0.05, \* p<0.1

Table 6 shows that *techdiff* is also related to relevant accounting variables. Columns 1 and 2 show that a one standard deviation higher *techdiff* is associated with 3.3% higher return on assets (ROA). Columns 3 and 4 show that *techdiff* is also positively associated with investment (capex over total assets) although the relation is more modest, with a one standard deviation higher *techdiff* associated with a capex to book assets ratio that is 0.5% higher relative to the firm's sample average.

## 4 Model

If R&D immediately converts into new products then firms' technological and product market positions are always identical. There is reason to think that the real optionality in firms' asset holdings is material, however. In this section, I develop a model that characterizes the effects that growth options on differentiated products have on the cross section of expected returns.

Consider a single real option to develop a new product. At any time the firm can pay  $I$  to embody the product, at which time it starts producing a continuous cash flow  $CF_t$ . The (potential) cash flow evolves according to:

$$\frac{dCF}{CF} = \nu dt + \beta dZ_M + \eta dZ$$

The CAPM holds<sup>3</sup>, and the market portfolio  $M$  follows:

$$\frac{dM}{M} = \mu dt + \zeta dZ_M$$

The risk free rate is fixed and for simplicity I set it equal to zero.

$D_t = \frac{CF_t}{\beta\mu - \nu}$  is the NPV of the embodied cash flows at time  $t$ , which has the same dynamics as  $CF_t$ . Consider the option to invest as a perpetual call option on  $D_t$ .  $D_t$  has risk neutral dynamics of

$$\frac{dD_t}{D_t} = (\nu - \beta\mu) dt + \beta\zeta dZ_M^{\mathbb{Q}} + \eta dZ = (\nu - \beta\mu) dt + \sigma d\hat{Z}^{\mathbb{Q}}$$

---

<sup>3</sup>The CAPM is only used for its simplicity. The effect of growth options on expected returns relative to assets in place is the same for any systematic risk factor in e.g. a multi-factor model.

The value of the option to invest  $V^D$  is:

$$V(D_t) = \sup_{\tau \geq t} \mathbb{E}_t^{\mathbb{Q}} \left[ (D_\tau - I)^+ \right]$$

which has corresponding PDE in the continuation region:

$$0 = \frac{1}{2} \sigma^2 D_t^2 V_{DD} + (\nu - \beta \mu) D_t V_D$$

with value-matching and smooth-pasting conditions:

$$V(D^*) = D^* - I$$

$$V_D(D^*) = 1$$

The firm optimally follows a rule where it invests as soon as  $D$  is above a threshold  $D^*$ :

$$D^* = \frac{B}{B-1} I, \quad B = 1 + 2 \frac{\beta \mu - \nu}{\sigma^2}$$

$$V(D_t) = (D^* - I) \left( \frac{D_t}{D^*} \right)^B$$

Assume that  $B > 1$  (otherwise there is no opportunity cost of waiting, so the firm never exercises the option).  $D^* > I$ , that is, firms optimally wait to invest even when the NPV is positive.  $D^*$  is decreasing in  $\beta$ , and for typical values of  $B$  this effect can be large. Hence, firms with low-beta projects wait longer to implement them. The beta of a particular growth option is

$$\hat{\beta} = \frac{\partial V}{\partial D} \frac{D}{V} \beta = B\beta$$

Thus, optionality increases the heterogeneity in expected returns of the underlying assets.

These are all standard results dating back to [McDonald and Siegel \(1986\)](#). Note that firms' behavior is quite different from that implied by models of innovation that do not feature optionality, such as quality ladder models ([Klette and Kortum \(2004\)](#); [Aghion et al. \(2005\)](#); [Akçigit and Kerr \(2010\)](#)) or new-product models ([Romer \(1990\)](#); [Kung and Schmid \(2013\)](#)). In those models, the NPV rule holds, firms invest immediately and there are no growth options.

Firms are collections of products and options on products. Each firm contains a set of  $J^E$  embodied products and  $J^P$  patents. The firm's book and market values equal

$$BOOK_t = \sum_{j \in J^E} I_j$$

$$MKT_t = \sum_{j \in J^E} D_t^j + \sum_{j \in J^P} V_t^j$$

For simplicity I assume that products differ only in their market betas. More differentiated products have a lower beta for two reasons; first, they face lower demand elasticity and have a higher pass through, and second, they are more profitable and hence face lower operating leverage ([Carlson, Fisher and Giammarino \(2004\)](#)).



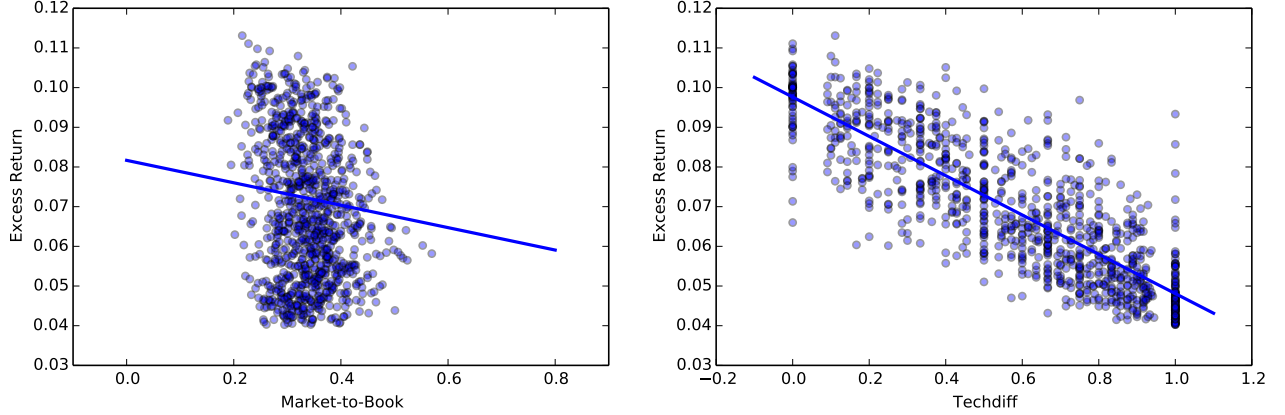


Figure 5: The figure plots the cross-section of expected returns of 1000 simulated firms with heterogeneous innovations.

Figure 5 plots the cross-section of expected returns that results from simulating 1000 firms in a calibrated economy with two types of products: nondifferentiated (Type 1:  $\beta = 1.5$ ) and differentiated (Type 2:  $\beta = 0.5$ ). The vertical axis indexes average excess returns, while the horizontal axes index market-to-book ratios and *techdiff* (here the firm's ratio of type-2 patents to type-1 patents). Appendix C describes the simulation procedure.

Panel A shows that the relation between market-to-book and simulated returns is negative. This is the mechanical value premium that Berk (1995) emphasizes. Firms with more differentiated products are less risky, so all else equal they have higher valuations and lower returns. The relationship between market-to-book and expected returns captures differences in returns due to firms' assets in place – that is, their product market position. Panel B shows that the relationship of simulated returns to *techdiff* is also negative, and is enormously stronger than that of Panel A because of the option leverage effect. In a joint regression using the simulated data, the coefficients of returns on market-to-book and *techdiff* are -0.028 and -0.048 with *t*-statistics of 5.6 and 45.7 respectively.

To sum up, the message of the model is that growth options amplify heterogeneity in the

expected returns of firm assets, and the resulting pattern in expected returns is not reflected in contemporaneous accounting measures. The predictions of the model are that:

1. Technological differentiation predicts lower returns after conditioning on accounting variables
2. Product market differentiation does *not* predict returns after conditioning on accounting variables
3. The effects of technological differentiation on returns are concentrated in growth firms

## 5 Returns

Monthly stock returns are from CRSP. I retain only U.S. firms traded on the NYSE, NASDAQ and AMEX from 1961 to 2008. I follow standard practice and update firms' accounting data in July using data from the fiscal year ending in the previous December at the latest. I drop any firms with a market value or book value of assets that is missing or less than \$10M. For the results in this section I index patents by their grant year, not their application year, and firms first appear in the panel the year after their first patent is granted.

Table 7 presents the results of Fama-MacBeth regressions of monthly excess returns on *techdiff* along with other variables and standard benchmarks. The first column shows that *techdiff* is negatively associated with subsequent returns after controlling for firm size (market cap), market-to-book ratio, one month reversal and 12 month momentum. The second column shows that the relationship persists after controlling for (citation-weighted) patent counts, R&D intensity and R&D productivity which have been previously shown to be associated with excess returns. These results are consistent with prediction #1 of the model.

The third column of Table 7 shows that measures of product market differentiation do not predict returns. These results are consistent with prediction #2 of the model, that the risk of assets in place is well measured by accounting variables.

Table 7: Results of Fama-MacBeth regressions of monthly excess returns for patenting firms from 1961-2008 on text-based technological differentiation (*techdiff*), measures of product market differentiation, and other variables. The dependent variable is annualized monthly excess returns,  $(rtn_{t+1} - r_t^f) \times 12$ . All independent variables have been standardized to have zero mean and unit variance.

	(1) $rtn - rf$	(2) $rtn - rf$	(3) $rtn - rf$
<i>techdiff</i>	-0.032*** (0.0068)	-0.034*** (0.0067)	-0.061*** (0.016)
<i>log(cites)</i>		0.020*** (0.0068)	
<i>RD/sales</i>		0.15 (0.13)	-0.0100 (0.015)
<i>HHI<sub>SIC</sub></i>		-0.0033 (0.0043)	-0.0054 (0.011)
<i>TNIC3tsimm</i>			0.017 (0.023)
<i>log(Size)</i>	-0.063*** (0.013)	-0.064*** (0.013)	-0.087*** (0.029)
<i>Book/Mkt</i>	0.018** (0.0081)	0.017** (0.0085)	0.000060 (0.018)
<i>rtn<sub>t-1</sub></i>	-0.12*** (0.0087)	-0.12*** (0.0094)	-0.080*** (0.016)
<i>rtn<sub>t-12,t-1</sub></i>	0.036*** (0.011)	0.040*** (0.011)	-0.018 (0.022)
Observations	496,741	406,698	156,266
R-squared	0.060	0.085	0.075

Standard errors in parentheses  
\*\*\* p<0.01, \*\* p<0.05, \* p<0.1

Table 8 shows that the association of greater technological differentiation with lower

returns is concentrated in growth firms, but is not strongly concentrated in small versus large firms.

Table 8: Results of Fama-MacBeth regressions of monthly excess returns for patenting firms from 1961-2008 on text-based technological differentiation (*techdiff*) interacted with market-to-book and size. The dependent variable is annualized monthly excess returns,  $(rtn_{t+1} - r_t^f) \times 12$ . All independent variables have been standardized to have zero mean and unit variance.

	(1) $rtn - rf$	(2) $rtn - rf$
<i>techdiff</i>	-0.053*** (0.010)	-0.0035** (0.0015)
<i>techdiff</i> $\times$ <i>Book/Mkt</i> <sub>quintile</sub>	0.0070*** (0.0026)	
<i>Book/Mkt</i> <sub>quintile</sub>	0.0057 (0.0056)	
<i>techdiff</i> $\times$ <i>Size</i> <sub>quintile</sub>		0.00044 (0.00032)
<i>Size</i> <sub>quintile</sub>		-0.0029*** (0.00072)
$\log(\text{Size})$	-0.065*** (0.013)	
<i>Book/Mkt</i>		0.0019*** (0.00065)
$rtn_{t-1}$	-0.12*** (0.0086)	-0.0098*** (0.00073)
$rtn_{t-12,t-1}$	0.037*** (0.011)	0.0030*** (0.00089)
Observations	497,311	496,741
R-squared	0.064	0.062

Standard errors in parentheses

\*\*\* p<0.01, \*\* p<0.05, \* p<0.1

Table 9: Panel A shows monthly portfolio returns after double independent sorts on size (market value of equity) and *techdiff* . Panel B shows portfolio betas against the *techdiff* HML factor.

Panel A: Excess Returns by Portfolio							
		<i>techdiff</i> Portfolio					
		1	2	3	4	5	HML
Size Portfolio	1	0.0129 (0.0070)	0.0079 (0.0044)	0.0095 (0.0041)	0.0071 (0.0035)	0.0047 (0.0033)	-0.0082 (0.0054)
	2	0.0090 (0.0048)	0.0026 (0.0036)	0.0003 (0.0033)	-0.0011 (0.0032)	-0.0026 (0.0031)	-0.0116*** (0.0034)
	3	0.0038 (0.0036)	0.0001 (0.0033)	-0.0024 (0.0028)	-0.0027 (0.0029)	-0.0035 (0.0029)	-0.0073*** (0.0023)
	4	-0.0001 (0.0028)	-0.0027 (0.0027)	-0.0032 (0.0027)	-0.0032 (0.0028)	-0.0039 (0.0026)	-0.0038** (0.0016)
	5	-0.0041 (0.0022)	-0.0047 (0.0024)	-0.0051 (0.0022)	-0.0042 (0.0023)	-0.0049 (0.0023)	-0.0008 (0.0014)

Panel B: Betas against HML						
		<i>techdiff</i> Portfolio				
		1	2	3	4	5
Size Portfolio	1	-2.56 (0.10)	-0.96 (0.08)	-0.90 (0.08)	-0.60 (0.07)	-0.34 (0.06)
	2	-1.56 (0.08)	-0.80 (0.07)	-0.65 (0.06)	-0.42 (0.06)	-0.32 (0.06)
	3	-1.01 (0.06)	-0.68 (0.06)	-0.38 (0.05)	-0.38 (0.05)	-0.16 (0.06)
	4	-0.51 (0.05)	-0.36 (0.05)	-0.25 (0.05)	-0.19 (0.05)	-0.10 (0.05)
	5	-0.25 (0.04)	-0.27 (0.04)	-0.12 (0.04)	-0.13 (0.05)	0.04 (0.04)

Table 9 shows the results of monthly independent double sorts on size (market capitalization) and *techdiff*. Panel A shows sizeable differences in monthly returns between the highest and lowest *techdiff* quintile, up to 1.16% per month in the second-smallest size quintile. The return difference is primarily among the smallest firms by market cap, and is absent in the largest size quintile.

Panel B shows betas of each of the portfolios against the monthly returns of the high-

minus-low *techdiff* portfolio. We see that there is a nearly monotonic relationship between the betas and *techdiff* quintiles, which is generally consistent with the “risk” story in this paper and less consistent with a “mispricing” story.

## 6 Conclusion

This paper uses patent text to impute firms’ technological positions. Firms’ location in technology space is related to but reliably different from their location in product market space per Hoberg and Phillips. I find that firms’ degree of technological differentiation predicts future product differentiation, but the reverse is not true.

Technological differentiation is strongly positively correlated with both average  $Q$  and  $TFP$ . Firms that are more technologically differentiated have lower returns conditional on standard benchmarks, while product market differentiation does not predict returns. The pattern is stronger in small growth firms and weaker in large value firms, and obtains both in portfolio sorts and in Fama-MacBeth regressions controlling for other known return predictors. These results are consistent with a simple model of real options on heterogeneous assets.

## References

- Aghion, Philippe et al., ‘Competition and Innovation: an Inverted-U Relationship’, *The Quarterly Journal of Economics*, 120 (2005):2, pp.701–728 ⟨URL: <http://qje.oxfordjournals.org/content/120/2/701.short>⟩.
- Akcigit, Ufuk and Kerr, William R., *Growth through heterogeneous innovations*, (National Bureau of Economic Research, 2010) – Technical report ⟨URL: <http://www.nber.org/papers/w16443>⟩.
- Alexopoulos, Michelle, ‘Read all about it! What happens following a technology shock?’ *The American Economic Review*, 101 (2011):4, pp.1144–1179 ⟨URL: <http://www.ingentaconnect.com/content/aea/aer/2011/00000101/00000004/art00005>⟩.
- Berk, Jonathan B, ‘A Critique of Size Related Anomalies’, *Review of Financial Studies*, 8 Summer (1995):2, pp. 275–286.
- Berk, Jonathan B., Green, Richard C. and Naik, Vasant, ‘Valuation and return dynamics of new ventures’, *Review of Financial Studies*, 17 (2004):1, pp.1–35 ⟨URL: <http://rfs.oxfordjournals.org/content/17/1/1.short>⟩.
- Carlson, Murray, Fisher, Adlai and Giammarino, Ron, ‘Corporate Investment and Asset Price Dynamics: Implications for the Cross-section of Returns’, *The Journal of Finance*, 59 (2004):6, pp.2577–2603 ⟨URL: <http://onlinelibrary.wiley.com/doi/10.1111/j.1540-6261.2004.00709.x/full>⟩.
- Cohen, Lauren, Diether, Karl and Malloy, Christopher, ‘Misvaluing innovation’, *Review of Financial Studies*, (2013), p.hhs183 ⟨URL: <http://rfs.oxfordjournals.org/content/early/2013/01/21/rfs.hhs183.short>⟩ – visited on 2014-10-30.

- Garleanu, Nicolae, Panageas, Stavros and Yu, Jianfeng, ‘Technological growth and asset pricing’, *The Journal of Finance*, 67 (2012):4, pp.1265–1292 ⟨URL: <http://onlinelibrary.wiley.com/doi/10.1111/j.1540-6261.2012.01747.x/full>⟩.
- Griliches, Z., ‘PATENT STATISTICS AS ECONOMIC INDICATORS: A SURVEY.’ *Journal of Economic Literature*, 28 (1990):4, pp. 1661–1707 ⟨URL: <http://elibrary.ru/item.asp?id=1561055>⟩.
- Griliches, Zvi, ‘Market value, R&D, and patents’, *Economics letters*, 7 (1981):2, pp.183–187 ⟨URL: <http://www.sciencedirect.com/science/article/pii/0165176587901145>⟩.
- Hall, Bronwyn H., Jaffe, Adam B. and Trajtenberg, Manuel, *The NBER patent citation data file: Lessons, insights and methodological tools*, (National Bureau of Economic Research, 2001) – Technical report ⟨URL: <http://www.nber.org/papers/w8498>⟩.
- Hirshleifer, David, Hsu, P. and Li, Dongmei, ‘Don’t Hide Your Light Under a Bushel: Innovative Diversity and Stock Returns’, *Available at SSRN 2117516*, (2012) ⟨URL: [http://rady.ucsd.edu/docs/ID\\_Returns\\_SSRN\\_submission.pdf](http://rady.ucsd.edu/docs/ID_Returns_SSRN_submission.pdf)⟩.
- Hirshleifer, David, Hsu, Po-Hsuan and Li, Dongmei, ‘Innovative efficiency and stock returns’, *Journal of Financial Economics*, 107 (2013):3, pp.632–654 ⟨URL: <http://www.sciencedirect.com/science/article/pii/S0304405X12001961>⟩.
- Hoberg, G. and Phillips, G., ‘Product market synergies and competition in mergers and acquisitions: A text-based analysis’, *Review of Financial Studies*, 23 (2010):10, pp. 3773–3811 ⟨URL: <http://rfs.oxfordjournals.org/content/23/10/3773.short>⟩.



- Hoberg, Gerard and Phillips, Gordon, ‘The Stock Market, Product Uniqueness, and Comovement of Peer Firms’, *Product Uniqueness, and Comovement of Peer Firms* (October 26, 2012), (2012) ⟨URL: [http://papers.ssrn.com/sol3/papers.cfm?abstract\\_id=2160846](http://papers.ssrn.com/sol3/papers.cfm?abstract_id=2160846)⟩.
- Imrohoroglu, Ayse and Tuzel, Selale, ‘Firm-Level Productivity, Risk, and Return’, *Management Science*, (2014) ⟨URL: <http://pubsonline.informs.org/doi/abs/10.1287/mnsc.2013.1852>⟩.
- Klette, Tor Jakob and Kortum, Samuel, ‘Innovating Firms and Aggregate Innovation’, *Journal of Political Economy*, 112 (2004):5, pp. 986–1018 ⟨URL: <http://www.jstor.org/stable/10.1086/422563>⟩.
- Kogan, L. et al., ‘Technological Innovation, Resource Allocation and Growth’, *Working Paper*, (2014) ⟨URL: [http://papers.ssrn.com/sol3/papers.cfm?abstract\\_id=2193068](http://papers.ssrn.com/sol3/papers.cfm?abstract_id=2193068)⟩.
- Kung, Howard and Schmid, Lukas, ‘Innovation, growth and asset prices’, *Duke University*, (2013) ⟨URL: <https://faculty.fuqua.duke.edu/~ls111/GrowthMay.PDF>⟩.
- McDonald, Robert and Siegel, Daniel, ‘The Value of Waiting to Invest’, *The Quarterly Journal of Economics*, 101 November (1986):4, pp. 707–728 ⟨URL: <http://www.jstor.org.libproxy.usc.edu/stable/1884175>⟩, ISSN 0033–5533.
- Nicholas, Tom, ‘Does innovation cause stock market runups? Evidence from the great crash’, *The American Economic Review*, (2008), pp. 1370–1396 ⟨URL: <http://www.jstor.org/stable/29730126>⟩.

- Packalen, Mikko and Bhattacharya, Jay, *Words in Patents: Research Inputs and the Value of Innovativeness in Invention*, (National Bureau of Economic Research, 2012) – Technical report ⟨URL: <http://www.nber.org/papers/w18494>⟩.
- Pastor, Lubos and Veronesi, Pietro, ‘Stock valuation and learning about profitability’, *The Journal of Finance*, 58 (2003):5, pp. 1749–1790 ⟨URL: <http://onlinelibrary.wiley.com/doi/10.1111/1540-6261.00587/abstract>⟩.
- Romer, Paul M., ‘Endogenous Technological Change’, *Journal of Political Economy*, 98 (1990):5 pt 2 ⟨URL: <http://individual.utoronto.ca/zheli/A2.pdf>⟩.
- Solow, Robert M., ‘Technical progress, capital formation, and economic growth’, *American Economic Review*, 52 (1962):2, pp. 76–86 ⟨URL: <http://sites-final.uclouvain.be/econ/DW/DOCTORALWS2004/bruno/vintage/solow.pdf>⟩.

## A Selecting Ngrams

To select words and phrases that summarize patents into vectors, I do the following:

For each patent:

1. Select body text (background, description and claims sections)
2. Remove words that are only numbers or only contain ATCG (genetic code)
3. Extract all single words and all two and three word phrases ('ngrams')
4. Filter out standard stopwords such as articles, prepositions, and place names from a stopwords dictionary
5. Compare each word in each ngram to the Wordnet lexicon to see if it can be used as an English noun
6. Keep only ngrams that contain exclusively English nouns with no stopwords

Across the corpus:

1. Count the number of documents (patents) in which each ngram occurs
2. Select ngrams that appear in no fewer than 1000 and no more than 1% of all ngrams in the sample (i.e. no more than 62,238 patents)

This process yields the 'master wordlist' of 64,369 words and phrases. Each patent is then represented as a Boolean vector with a one in positions where the patent contains that word or phrase and a zero if not. Patents are dropped from the vectorized sample if they contain fewer than 1000 words or fewer than 10 ngrams.

## B Individual Patents

### B.1 Technology Measures by Individual Patent

I compute the cosine similarity of each patent vector with every other patent vector. This is a commonly used measure of two documents' similarity; for patent vectors  $i$  and  $j$  it equals the dot product

$$\cos_{ij} = \left( \frac{pvec_i}{\|pvec_i\|} \right) \cdot \left( \frac{pvec_j}{\|pvec_j\|} \right)$$

The first patent-level measure I calculate is patent *quality*. This measure is based on two sub-measures: *influence*, which is the median similarity of the focal patent to the 100 nearest neighbouring patents from subsequent years, and *originality*, which is one minus the median similarity of the focal patent to the 100 nearest neighbouring patents from previous years. Originality measures how uncommon a patent's ngrams – and the specific mix of ngrams in the patent – were in patents that precede it. Influence measures how common they were in patents that follow. Text-based *quality* is originality plus influence – or equivalently, the difference in similarity of patents that postdate the focal patent minus those that predate the focal patent. For example, ngrams relating to the Internet do not appear at all in the 1970s and 80s but appear increasingly often in the 1990s and 2000s. Thus, a patent related to the internet that appears in the 1990s will be generally scored as having high originality, influence, and quality.

Patents' originality and influence scores are strongly negatively correlated ( $\rho = -0.71$ ), because some ngrams are more common than others throughout the sample. A patent related to microprocessors or recombinant DNA tends to have both low originality and high influence, relative to a patent related to hockey helmets. The *quality* measure avoids this issue: if a patent's similarity to its decedents and antecedents is the sum of a patent-specific

component and a non time varying component specific to its area, then *quality* differences out the latter.

The second patent-level measure is *vintage*. For each ngram in the list, I note the first year that that word or phrase appears in any patent. *vintage* is defined as the dot product of a patent’s Boolean vector with the vector of first-mention years, divided by the number of entries in the vector:

$$vintage_i = \frac{pvec'_i \cdot firstmention}{pvec'_i \cdot \mathbf{1}}$$

*vintage* is thus the average year of first mention across all the ngrams contained in the patent. It reflects how current the language in the patent is, without direct reference to other patents.

## B.2 Validation: Patent Renewals

A granted patent is valid for up to twenty years; for all patents granted after December 1981, the patent expires at four, eight and twelve years after the grant date unless a renewal fee is paid. As of 2013 the fees were \$1600, \$3600 and \$7400. Thus, one measure of whether or not a patent’s value exceeds some minimal value is whether it was renewed or not. Of the 1.61 million patents in the sample granted from 1982 to 1998, 84%, 63% and 44% were renewed after four, eight and twelve years respectively.

Table 10 shows the results of regressing a dummy variable that equals one if the patent was renewed after twelve years on our measures of patent quality. Citations are strongly correlated with renewal. Some of this is likely a selection effect because both are ex post outcomes: conditional on patent value, a patent that is not renewed is less likely to be cited.

Text-based *quality* and *vintage* are potentially subject to a similar critique, because they are based on the list of informative ngrams which are selected from their usage throughout the entire sample. But because the selection effect is at the level of ngrams rather than the patent itself, the bias for the text-based measures is likely to be less than it is for citations. We see that both text-based measures are strongly positively correlated with subsequent renewal. A one standard deviation increase in *quality* or *vintage* is associated with an increase in the likelihood that the patent will be renewed after 12 years by 3.6% and 7.0%, respectively.

Table 10: Logit regressions of patent renewal on text based measures, for all U.S. utility patents granted from 1982 to 1998. The dependent variable is a dummy variable that equals one if the patent is renewed at the twelve year mark. The independent variables are standardized to unit variance, and the coefficients report the estimated marginal effects. The standard errors are clustered by year.

	Logit 12yr Renew	Logit 12yr Renew	Logit 12yr Renew
<i>quality<sub>i</sub></i>		0.036*** (0.0040)	0.014*** (0.0029)
<i>vintage<sub>i</sub></i>		0.070*** (0.0043)	0.062*** (0.0043)
<i>ln(cites<sub>i</sub>)</i>	0.099*** (0.0030)		0.084*** (0.0023)
Observations	1,617,938	1,617,938	1,617,938
Fixed Effects	Year	Year	Year
Pseudo R-squared	0.032	0.026	0.043
Robust standard errors in parentheses			
*** p<0.01, ** p<0.05, * p<0.1			

## C Simulation

To illustrate the effects of growth options on the cross section of expected returns, I simulate a cross section of 1000 firms. In the economy the CAPM holds with  $\mu = 0.06$ ,  $r_f = 0$ . There are two types of products. Type 1 products are nondifferentiated ( $\beta = 1.5$ ), Type 2 products are differentiated ( $\beta = 0.5$ ). Otherwise the products are identical with parameters  $I = 1$ ,  $\sigma = 0.2$ ,  $\nu = -0.02$ ,  $D_{initial} = 1$ . Firms have arrival probabilities of innovations that are drawn from a uniform joint distribution on  $[0, 0.1] \times [0, 0.1]$ . Each step consists of:

1. Cash flow shocks hit all existing products and options
2. Firms exercise options with  $D_t > D_j^*$
3. New innovations arrive
4. Firms update their book values (i.e. issue a quarterly report)

I simulate 1000 firms for 10,000 steps which is enough to reach the stationary distribution (results for 20,000 steps look the same).