

# NBER WORKING PAPER SERIES

## THE DIFFUSION OF NEW TECHNOLOGIES

Nicholas Bloom  
Tarek Alexander Hassan  
Aakash Kalyani  
Josh Lerner  
Ahmed Tahoun

Working Paper 28999  
<http://www.nber.org/papers/w28999>

NATIONAL BUREAU OF ECONOMIC RESEARCH  
1050 Massachusetts Avenue  
Cambridge, MA 02138  
November 2023

The views expressed herein are solely those of the authors and do not necessarily reflect those of the Federal Reserve Bank of St. Louis, the Federal Reserve System, or the National Bureau of Economic Research. We thank audiences at the Applied Machine Learning webinar, Atlanta Fed, Babson College, the Bank for International Settlements, Baruch College, Bocconi University, CKGSB, Columbia University, Dartmouth University, Duke University, Durham University, ETH, Georgia State University, Harvard University, the Korea-American Economic Association, the London Business School, the London School of Economics, Michigan State University, Northwestern University, Nova Business School, New York University, the Ohio State University, the Royal Bank of Australia, Stanford University, the Toulouse Network on Information Technology, the University of British Columbia, the University of California at Santa Barbara and Santa Clara, the University of Chicago, the University of Maryland, the University of Michigan, the University of Minnesota, the University of North Carolina, the University of Southern California, the University of Texas, the University of Washington, Yeshiva University, and the NBER 2021 Summer Institute and Fall 2021 EFG meetings and the 2022 Society for Economic Dynamics annual meeting for helpful comments. Special thanks go to Lisa Kahn for sharing data, Bledi Taska for help on BGT data queries, Ben Jones and Chad Syverson for excellent discussions, and William Hartog and Jared Simpson for excellent research assistance. We thank Scarlett Chen, Nick Short, Corinne Stephenson, and Michael Webb for assistance in conceptualizing and researching early versions of this project. Funding for this research was provided by Harvard Business School, the Institute for New Economic Thinking, the Kauffman Foundation, the Sloan Foundation, the Toulouse Network on Information Technology, and the Wheeler Institute. Bloom and Lerner have received compensation from advising institutional investors in venture capital funds, venture capital groups, and governments on venture capital topics. All errors and omissions are our own.

NBER working papers are circulated for discussion and comment purposes. They have not been peer-reviewed or been subject to the review by the NBER Board of Directors that accompanies official NBER publications.

© 2021 by Nicholas Bloom, Tarek Alexander Hassan, Aakash Kalyani, Josh Lerner, and Ahmed Tahoun. All rights reserved. Short sections of text, not to exceed two paragraphs, may be quoted without explicit permission provided that full credit, including © notice, is given to the source.

The Diffusion of New Technologies

Nicholas Bloom, Tarek Alexander Hassan, Aakash Kalyani, Josh Lerner, and Ahmed Tahoun  
NBER Working Paper No. 28999

July 2021

JEL No. O31,O32

### **ABSTRACT**

We identify phrases associated with novel technologies using textual analysis of patents, job postings, and earnings calls, enabling us to identify four stylized facts on the diffusion of jobs relating to new technologies. First, the development of new technologies is geographically highly concentrated, more so even than overall patenting: 56% of the economically most impactful technologies come from just two U.S. locations, Silicon Valley and the Northeast Corridor. Second, as the technologies mature and the number of related jobs grows, hiring spreads geographically. But this process is very slow, taking around 50 years to disperse fully. Third, while initial hiring in new technologies is highly skill biased, over time the mean skill level in new positions declines, drawing in an increasing number of lower-skilled workers. Finally, the geographic spread of hiring is slowest for higher-skilled positions, with the locations where new technologies were pioneered remaining the focus for the technology's high-skill jobs for decades.

Nicholas Bloom  
Stanford University  
Department of Economics  
579 Jane Stanford Way  
Stanford, CA 94305-6072  
and NBER  
nbloom@stanford.edu

Josh Lerner  
Harvard Business School  
Rock Center 214  
Soldiers Field  
Boston, MA 02163  
and NBER  
jlerner@hbs.edu

Tarek Alexander Hassan  
Department of Economics  
Boston University  
270 Bay State Road  
Boston, MA 02215  
and NBER  
thassan@bu.edu

Ahmed Tahoun  
London Business School  
26 Sussex plc, Regent's Park  
London NW1 4SA  
United Kingdom  
atahoun@london.edu

Aakash Kalyani  
Department of Economics  
Boston University  
270 Bay State Road  
Boston, MA 02215  
aakashk@bu.edu

The dataset constructed as part of this paper, as well as relevant code. is available at  
[www.techdiffusion.net](http://www.techdiffusion.net)

## 1. Introduction

Economists have long recognized that the development of novel technologies is inexorably linked to economic growth. Many studies have sought to understand whether the benefits from adopting new technologies accrue primarily to inventors, early investors, highly skilled users, or to society more widely through, for instance, employment and income growth.<sup>2</sup> Substantial concerns remain, however, about as to the implications of new technologies, including whether they contribute to income inequality (e.g., do technology-enabled jobs spread beyond college graduates?) and regional inequality (do technology jobs spread outside Silicon Valley?).<sup>3</sup>

One key obstacle to resolving these questions is that it has proven difficult to measure the development and spread of multiple technological advances in a single framework and to systematically identify those innovations that affect jobs and businesses. In this paper, we use the full text of millions of patents and job postings and hundreds of thousands of earnings conference calls to make progress on this issue. We develop a flexible methodology that allows us to determine which (sets of) technological innovations most affected businesses over the past two decades, trace these back to the locations and firms where they emerged, and track their diffusion through regions, occupations, and industries over time. We then use our newly created data to establish key stylized facts about the development and diffusion of new technologies across space and skill levels.

The first step of our analysis is to develop a methodology for systematically identifying one, two, and three-word phrases (unigrams, bigrams, and trigrams) associated with new technologies through a series of systematic rules, whose robustness we verify through various diagnostic tests. To this end, we intersect information from multiple large corpora of text. First, we use the full text of U.S. patents with application years between 1976 and 2014 to isolate phrases that appear in multiple patents but did not exist before 1970. That is, we isolate new language specific to influential innovations made in the past 40 years. Second, we search for these phrases in Wikipedia to identify which of these new phrases are primarily associated with pages describing new technologies, as opposed to newly recognized problems (such as “climate change”) or new management terms (such as “performance metrics”). This procedure identifies 1,899 new

---

<sup>2</sup> See, for example, Katz and Murphy (1992), Krusell et al. (2000), Piketty and Saez (2003), Autor et al. (2008), Goldin and Katz (2009), Acemoglu and Autor (2011), and Song et al. (2019).

<sup>3</sup> See Tyson and Spence (2017) and Vance (2022) for popular articulations of such concerns.

technology phrases, which we can group into 1,286 unique Wikipedia pages describing new technologies. We refer to these groups of new technology phrases as “technologies.”

After establishing our list of new technologies, we then identify patents and job postings that mention these technologies. We use patent inventor addresses to identify the locations where each of the technologies was developed and patent application years to pinpoint the year in which the technology experienced the first large acceleration in patent references (its “emergence year”). We then cross-reference our list of technology phrases with the full text of online job postings to identify 51 million jobs advertised between 2010 and 2019 that mention these technologies. These granular data uniquely allow us to track the spread of new technologies along a dimension of crucial importance to policymakers: jobs. In particular, we examine the evolution of the number, location, and skill requirements of job postings associated with these new technologies.

In a final step, we use the full text of earnings conference calls held by listed firms between 2002 and 2019 to flag those new technologies that disrupted business during this period, in the sense that they are referenced in at least 100 of these important conversations between firm executives and investors. This last step yields 276 “disruptive” new technologies (unique Wikipedia pages). The most prominent of these new disruptive technologies include “cloud computing,” “smart phone,” and “machine learning.”

Strikingly, these few disruptive technologies account for the lion’s share of the variation in our data: on average, each disruptive technology is mentioned in 141,633 job postings between 2010 and 2019, whereas new technologies mentioned in less than ten earnings calls average only 1,165. In total, our disruptive technologies are mentioned in 39 million job postings (or about 77 percent of all job postings mentioning any new technology) and 33.1% of patents granted by the U.S. Patent and Trademark Office (USPTO) to domestic inventors with application years between 1976 and 2014. In this sense, the few innovations that feature prominently in managers’ discussions are also the ones that are prominent in job postings and associated with the most widespread patenting activity. We therefore pay special attention to these disruptive technologies throughout our analysis.

Our key results are as follows.

First, the locations where new technologies are developed are geographically highly concentrated. This concentration is particularly pronounced for disruptive technologies: 32.5% of patents

mentioning any new technology and 42.1% of patents mentioning a new disruptive technology emerge from just five urban areas: San Jose, San Francisco, New York, Seattle, and Boston.

Based on early patenting activity around the time of each technology's emergence year, we identify which urban areas housed the majority of early patenting for each of our new technologies. We term these urban areas "pioneer locations."

Again, these pioneer locations for new technologies are highly concentrated, particularly so for disruptive technologies, those with the biggest impact on businesses and job postings. Collectively, 56.3% of the economically most impactful technologies come from just two U.S. locations, Silicon Valley and the Northeast Corridor. (Locations in California collectively host a remarkable 41.0% of pioneer locations of disruptive technologies.) This extreme concentration is particularly important because new technologies alter the composition of the local job postings in their pioneer locations for long periods of time, as we show below.

Second, despite this highly skewed initial distribution of pioneer locations, as technologies mature and the number of new jobs related to them grows, they gradually spread geographically. Our favored measure of geographic concentration, the coefficient of variation of the share of jobs associated with a new technology across the 917 core-based statistical areas (CBSAs) in the United States, drops by 18.5% in the first decade after its emergence. Nevertheless, the implied years to full dispersion across CBSAs is in excess of 50 years, well beyond the horizon of many individuals or policymakers.

Third, while initial hiring is heavily biased towards high-skilled jobs, the mean required skill level of the jobs associated with new technologies declines over time, reflecting a broadening of the types of jobs that adopt a given technology. Specifically, we estimate that, in the year of the average technology's emergence, 57.1% of the initial jobs relating to this new technology require a college degree – a substantial skill bias relative to 30.3% of respondents in the 2015 ACS<sup>4</sup> that hold one. This gap declines by 0.23 percentage points per year, so that 30 years after a technology's emergence year, on average about 50.2% of job postings relating to it still require a college degree.

---

<sup>4</sup> In the 2021 ACS, this proportion is somewhat higher, 36.4%.

Fourth, low-skill jobs associated with a given technology spread out across space significantly faster than high-skill jobs, which tend to remain concentrated for long periods of time within the pioneer locations that originally developed the technology.

A key implication of these patterns is that new technologies appear to yield long-lasting benefits for the pioneer locations where they were originally developed. These locations host a disproportionate share of high-skilled jobs relating to these new technologies for about four decades after their year of emergence. In short, this concentration of disruptive innovation in a handful of urban centers engenders large and persistent regional disparities in economic opportunity, giving a handful of U.S. locations a lasting advantage in high-skill job postings.

To shed light on the mechanisms underlying these patterns, we study the context in which the new technology is mentioned within a given job posting. We find that much of the regional spread of new technologies is driven by low-skill jobs associated with their *use*, whereas jobs related to their research, development, and production (RDP) remain persistently concentrated in and around their pioneer locations. That is, our evidence suggests that pioneer locations that initially developed a technology retain a long-term advantage, because they retain the technology's RDP for long periods of time.

We also show evidence to suggest that some of the observed skill broadening of new technologies is driven by standardization that allows for the use of these technologies by lower-skill workers as the technology matures. By contrast, training and experience with the new technology do not appear to be major drivers of this process.

We conduct a large number of robustness checks, replicating our main results using a wide range of different variations. For example, we repeat our analysis with phrases of different lengths (such as unigrams and trigrams), conduct a human audit of technology phrases and technologies, and use alternative methods for pinpointing pioneer locations and emergence years. Throughout all of these variations, our main findings remain unchanged.

We note three main caveats to our interpretation. First, all of our results regarding jobs rely on the analysis of job postings. In this sense, they measure the characteristics of open positions, but not necessarily the characteristics of the jobs that get filled. Second, by its very nature, our data speak to job openings relating to novel technologies but not to the possible destruction of existing positions by these technologies. Finally, a third concern is what Merton (1968) termed

“obliteration by incorporation”: when a technology becomes so widely diffused that it is no longer mentioned specifically in job postings. For this reason, we focus on relatively recent technologies, rather than ones, such as electricity or air-conditioning, that have been around for so long that they became normalized.

Our work builds on a large literature that studies the relationship between technology and labor markets. One strand of this literature studies the diffusion of technology. This literature has focused on patterns in a single specific (though important) new technology, from computers (Autor et al., 2003) to broadband (Akerman et al., 2015) to robots (Acemoglu and Restrepo, 2020) to artificial intelligence (Agrawal et al., 2019; Webb, 2020). Other studies have focused on specific innovations during important historical episodes (Griliches, 1957; Goldin and Katz, 1998; Squicciarini and Voigtlander, 2015; Caprettini and Voth, 2020).<sup>5</sup> Comin and Hobijn (2004, 2010) characterize the diffusion of 15 technologies across 166 countries, employing a variety of measures of technological utilization at the country level.<sup>6</sup> We contribute to this literature by identifying hundreds of new technologies, pinpointing their geographic origins, and tracking their spreads across job postings, skill levels, and geographies within the United States.

A second strand is the literature on technology and inequality. Many of these works have sought to estimate the skill bias of technical progress (e.g., Katz and Murphy, 1992; Krueger, 1993; Berman et al., 1994; Autor et al., 1998; Goldin and Katz, 2008; Autor et al., 2008; Michaels et al., 2014; Song et al., 2019). The near-universal approach in this literature is to infer an increase in the demand for skilled labor over time from changes in observed wage differentials – in effect, documenting a change in the economy’s aggregate production function. Our work complements this literature by observing this skill bias of technical progress directly: for instance, 57.1% of early jobs involved with new technologies require a college degree.<sup>7</sup>

Closely related, Caselli (1999), Acemoglu et al. (2012), and Acemoglu and Restrepo (2018) study theoretically the forces that drive automation, the substitution of capital for labor, and inequality.

---

<sup>5</sup> Recent work has examined the importance of supply and demand factors for the speed of diffusion (e.g., Popp, 2002; Acemoglu and Linn, 2004; Greenstone et al., 2010; Moser et al., 2014; Moscona, 2020; Arora et al., 2021). Mokyr (1992) and Gordon (2016) trace out the impact on economic development of a range of great inventions.

<sup>6</sup> A large, related literature studies the role of trade and multinational production in facilitating the diffusion of technology. Recent examples include Buera and Oberfield (2020) and Lind and Ramondo (2022).

<sup>7</sup> Notably, Goldin and Katz (1998) show that the introduction of new manufacturing processes during the early 19<sup>th</sup> century increased the demand for skilled labor. Krueger (1993) shows that workers who use computers at work earn higher wages.

A key result in this literature is that balance in the “race between man and machine” arises endogenously if the use of new technologies spreads from high-skill to low-skill occupations over time. We contribute by providing the first direct evidence that this skill broadening indeed occurs systematically for a broad range of technologies. In addition, this theoretical literature argues that one mechanism underpinning skill broadening is that new technologies evolve over time into a standardized form that can more readily be used by less educated workers. We provide empirical evidence supporting this standardization mechanism.

A third broad literature examines clustering in entrepreneurial activity and innovation. A number of papers have highlighted persistent advantages in entrepreneurship (Glaeser et al., 2015) and innovation (Moretti, 2021) that certain urban areas enjoy and highlighted mechanisms such as employee mobility across new ventures (Gompers et al., 2005) and localized knowledge spillovers (e.g., Jaffe et al., 1993). We contribute to this literature by providing a systematic approach to identifying and studying pioneer locations. We characterize their distribution across the United States and show there is a general relationship between successful innovation, early employment in a new technology, and the long-term advantage that these locations enjoy in high-skill employment.

Finally, our work adds to a growing literature in economics using text as data. A number of recent papers have used newspaper articles, patents, and firm-level communications to measure concepts that are otherwise hard to quantify (e.g., Hoberg and Phillips 2016; Baker et al., 2016; Hassan et al., 2019, 2021; Bybee et al., 2020; Handley and Li, 2020; Flynn and Sastry, 2020; Kelly et al., 2021; and Sautner et al., 2023). We focus primarily on the full text of job postings, which has received relatively less attention.<sup>8</sup> An important paper by Kogan et al. (2022) intersects information from patents and task descriptions to measure the direction of technical progress. Our work adds to this literature by introducing a flexible methodology for analyzing the origin and spread of innovations by intersecting multiple large corpuses of texts.

The remainder of this paper is structured as follows. In Section 2, we discuss how we identify and characterize new technologies in the data. Section 3 studies the spatial concentration of the development of new technologies. In Section 4, we explore the diffusion of activity across regions

---

<sup>8</sup> A notable exception is the work by Abis and Veldkamp (2020), who use job descriptions to identify financial analysis positions that leverage machine learning.



and the associated mechanisms. We present our analysis of the skill-broadening results in Section 5. Section 6 examines diffusion across occupations, industries, and firms. Section 7 presents robustness checks. The final section concludes the paper.

## 2. Identifying and Characterizing Technological Innovations

Our first objective is to identify a list of phrases describing *influential technological innovations developed since 1976*.

We use the term *technological innovation* in the sense of Schmookler (1966) and Jewkes et al. (1969), who distinguish technological from scientific innovation – the former being a set of specific and applied techniques, products, and processes (our focus here) while the latter is a set of general principles. This motivates our use of patents (as opposed to scientific research papers) as a text source.<sup>9</sup> We further distinguish technological from managerial knowledge. While Syverson (2011) and Bloom et al. (2016) argue that managerial rather than technological knowledge can account for substantial differences in total factor productivity across firms, we deliberately focus *on technological but not managerial knowledge* when we require that the new language we isolate from patents describe technologies.<sup>10</sup> By *influential* and *developed since 1976* we mean those innovations mentioned repeatedly in highly cited patents and those that went through a major acceleration in patenting activity after 1976.

We now describe in more detail how we operationalize these concepts in the data.

### a. Step 1: Identify phrases associated with influential innovations

We begin by examining patent filings with the USPTO. By law, patents must describe their technological innovation and (at least some) key ways in which it is applied.<sup>11</sup> Because of the

---

<sup>9</sup> The U.S. patentability standard requires an invention not to be obvious “to a person having ordinary skill in the art” (35 U.S.C. 103), but not an abstract idea, law of nature, or natural phenomenon (35 U.S.C. 101). See the discussions, for example, by the Supreme Court in *Alice Corp. v. CLS Bank International*, 573 U.S. 208 (2014) at 216 and *Mayo Collaborative Servs. v. Prometheus Labs., Inc.*, 566 U.S. 66 (2012) at 71.

<sup>10</sup> The OECD’s Oslo Manual (2005) elaborates on this distinction, providing many examples of what would and would not be included in the two categories.

<sup>11</sup> This requirement is stipulated in the legal concept of “reduction to practice,” 35 U.S.C. 112(a).

importance of the U.S. market, inventors worldwide typically file important discoveries with the USPTO.<sup>12</sup>

We collect all utility patents awarded to U.S. inventors with application years between 1976 and 2014, a total of approximately three million patents. We focus not just on the front page of the award, which has been the focus of much of the earlier analytic literature, but on the entire text of these patents. Representative parts of a patent are reproduced in Appendix Figure 1. For more details on this collection process refer to Section 1.1 of the Data Appendix.

To reduce the dimensionality of this voluminous body of text, we remove stop words (such as “of,” “the,” and “from”) following Kelly et al. (2021) and Gentzkow et al. (2019) and represent each patent’s remaining text by a vector of all two-word combinations (“bigrams”) that appear at least twice in the patent, leaving us with 17 million unique bigrams. In our main specification, we focus on bigrams because they are less ambiguous than single-word keywords. For example, while words like “autopilot” or “cloud” could have a variety of colloquial meanings, “autonomous vehicle” and “cloud computing” are much less ambiguous (e.g., Tan et al., 2002; Bekkerman and Allan, 2004). In Section 7 (robustness), we show that our results extend readily to including unigrams (one-word) and trigrams (three-word combinations), though unigrams generally appear to produce noisier results and trigrams add little to the analysis once bigrams are accounted for.

We next seek to isolate those bigrams that are novel and associated with influential innovations. First, we focus our attention on bigrams associated exclusively with *novel* innovations by dropping “non-novel” bigrams that were in common use before 1970. To this end, we select all text dating prior to 1970 from the Corpus of Historical American English, a representative sample of text constructed by linguists from prominent sources (Davies, 2009) that reflects everyday use of English up to 1970. We then remove any bigram appearing in this source (for instance, “equipment used”) from our list of bigrams obtained from patents, leaving us with 1.5 million exclusively “novel” bigrams.

---

<sup>12</sup> About half of all patent applications to the USPTO are filed by residents of foreign countries (USPTO, 2020). This pattern reflects the fact that patent protection in a given nation depends critically on having a patent issued in that specific nation. Important discoveries (the focus of our analysis) are therefore disproportionately likely to be filed in major patent offices worldwide (Lanjouw et al., 1998).

Second, to identify bigrams associated with *influential* innovations, we retain only those novel bigrams that appear in patents accumulating a total of at least 1,000 patent class and year-normalized citations.<sup>13, 14</sup> This leaves us with 36,563 novel and influential bigrams from patents.

b. Step 2: Identifying technological innovations using Wikipedia

A review of these novel and influential bigrams from patents suggests they fall into three broad categories. Some describe technological innovations, such as “fingerprint sensor,” “monoclonal antibody,” or “OLED display.” Others refer to new (or increasingly visible) problems, such as “greenhouse gases,” “Parkinson’s disease,” or “carbon footprint.” Yet others refer to areas that may have seen substantial new developments or management attention but are not new technologies, such as “account management,” “performance metrics,” and “business model.” (Appendix Table 1 shows examples.) As discussed above, we want to focus on bigrams in the first category, not the other two.

To isolate bigrams describing *technological* innovations, we employ Wikipedia entries. We first match each novel and influential bigram to a Wikipedia page by entering it into the Wikipedia search engine and selecting the highest-ranked entry if it mentions the bigram either in the title or the summary or it mentions the bigram at least 10 times in the body of the entry. Bigrams that do not meet these criteria (those without a Wikipedia page) are deleted.

The second step exploits the standardized nature of Wikipedia page entries. Entries describing technological innovations tend to feature sections containing the words *application(s)*, *use(s)*, *type(s)*, *operation*, *characteristic(s)*, *feature(s)*, *device(s)*, *technical*, and *commercial* in their titles. (Appendix Figure 2 provides examples of two Wikipedia pages with these features.) By contrast, pages dedicated to new problems or management innovations tend to feature sections and/or titles that contain the words *responses*, *mitigation*, *problems*, *causes*, *signs*, *symptoms*, *adverse effects*, *management*, *manager*, *risk assessment*, *business model*, *distribution model*, *customer*, *strategy*, and *service provider*. To focus on bigrams associated with technological innovations, we thus

---

<sup>13</sup> Following Lerner and Seru (2022), normalized citations for a patent  $p$  are calculated as:  $\frac{Citations_p}{Avg_{\tau,t}(Citations_{p'})}$ .

$Citations_p$  is the number of citations is received as of 2018 by a patent filed in four-digit Combined Patent Classification (CPC) technology class  $\tau$  in year  $t$ .  $Avg_{\tau,t}(Citations_{p'})$  is the average number of citations received by all patents filed in technology class  $\tau$  in year  $t$ .

<sup>14</sup> For computational reasons, it is necessary to restrict to a subset of the 1.5 million novel bigrams before cross-referencing with other corpuses (steps 2-4). However, where exactly we draw the boundary between influential and non-influential bigrams (1000 normalized citations) has little effect on our results, as discussed in Section 7.

retain only those that are matched with a Wikipedia page with at least one section from the former list, but none of the latter.

This algorithm returns a list of 4,277 bigrams associated with influential technological innovations, which we can conveniently group by the 2,746 unique primary Wikipedia pages that they are associated with. For ease of reference, we refer to these bigrams as “technology bigrams” and their groupings as “technologies,” which we label by the Wikipedia page’s title.<sup>15</sup> Appendix Table 1 provides examples of bigrams that passed and failed this Wikipedia filtering. For further details on scraping and processing Wikipedia pages, refer to Section 1.3 of the Data Appendix.

c. Step 3: Characterizing technologies using patents and earnings calls

To learn more about when and where each technology was developed, we next cross-reference our list of technology bigrams with our corpus of patents.<sup>16</sup> First, to obtain a measure for each technology’s age, we calculate for each bigram the first episode of accelerated patenting. In particular, we first calculate the number of cite-weighted patents (normalized as described in Section 2.a) mentioning the bigram filed in each calendar year. Due to the variability of the patent counts, we smooth the series by taking a centered five-year moving average. Finally, we mark the first year in which (a) the technology reaches 100 cite-weighted patents and (b) the next five years had at least 10% annual growth in the (smoothed) weighted patent filings. For ease of reference, we refer to this year as the bigram’s “emergence year.”

This process is illustrated in Figure 1, which depicts the time series and the emergence year for four technology bigrams. Digital video, for instance, emerges in 1986, as the time series grows by at least 10% for five consecutive years through 1991. Using this definition, we assign an emergence year after 1976 to 1,899 technology bigrams (1,286 technologies). The remaining bigrams exhibit no single five-year period of accelerated growth in our sample, and thus predominantly describe older technologies (such as diesel fuel and whey protein). In Section 7, we show our results are robust to using a range of other plausible approaches to defining each bigram’s

---

<sup>15</sup> We undertake a human audit to check each individual bigram and run robustness tests using this edited sample in Section 7. The bigrams deleted in the human audit are marked with asterisks in Appendix Table 3.

<sup>16</sup> When cross-referencing our technology bigrams with patents and other corpuses, we generally allow for all forms of the bigram, including singular, plural, and concatenations. We require the bigram to appear at least twice in the patent.

emergence year. The key is simply to obtain some meaningful distinction between older and newer innovations.

Second, to identify regions pioneering the early development of a technology, we identify the CBSAs that collectively account for a majority of early patents mentioning the technology. In particular, for each technology bigram, we calculate the number of patents in each CBSA within the first ten years of the bigram's emergence year. We then sort the CBSAs by the number of patents mentioning that bigram and denote those CBSAs with the most of these patents that collectively account for at least 50% of the total patents mentioning the technology bigram in this period as pioneer locations. Thus, if the top four CBSAs accounted for 25%, 15%, 12%, and 8% of the patents containing a bigram in this period, the first three would be coded as pioneer locations.

Third, we can gauge the extent to which a given technology poses economic challenges or opportunities to incumbent firms by cross-referencing our list of technologies with the full text of 321,373 corporate earnings calls held by 11,905 publicly held companies and compiled by Refinitiv EIKON between 2002 and 2019. Publicly traded firms hold quarterly earnings calls to discuss results and the companies' prospects. These calls (and the transcripts that we analyze) consist of a presentation by management (typically the chief executive and/or chief financial officer) and then questions posed by investors and analysts with answers by the executives. They have been shown to be indicators of some of the most important issues facing these organizations (Bushee et al., 2003; Matsumoto et al., 2011; Hassan et al., 2019, 2021).<sup>17</sup> To gauge the extent to which each technology features in the conversations at these large listed firms, we record the number of unique earnings calls in which each of our technologies is mentioned.

Table 1 gives a flavor of these data. It shows the top technology, as measured by the number of earnings calls mentioning it, by year of emergence of the technology, as well as its associated bigrams. Top technologies emerging in the late 1970s and early 1980s include the hard disk drive, barcode reader, and personal computer. The mobile phone emerges in 1985, followed by digital video and debit cards. The 1990s brought machine learning and the hybrid electric vehicle. The top technologies from the 2000s include the smartphone, social networking, and the self-driving car. Taken together, these technologies appear to accurately reflect the changing nature of technological innovation over the past decades. Appendix Table 3 lists all new technologies that

---

<sup>17</sup> Some examples of mentions of bigrams are contained Appendix Table 2.

are mentioned in more than 100 earnings calls. While we make no claim of completeness, we argue they constitute perhaps the most representative sample of economically influential technological innovations constructed to date.

Table 2 provides examples of the pioneer locations for a number of technology bigrams. For example, pioneer locations for machine learning (a technology that emerged in 1994 according to our measure) are New York, Seattle, San Jose, and San Francisco, whereas digital imaging's pioneers are Rochester (Kodak's headquarters), San Jose, San Francisco, and Fort Collins (the longtime home of Hewlett Packard's desktop and peripherals business).

Appendix Table 4 shows, for selected states, the technology where the state most dominated early innovation; that is, the technology where the state contributed the largest share of early patenting. The table again shows intuitive patterns. For example, Massachusetts accounts for 13.6% of the early patenting in the technology "antibody-drug conjugate," and similarly, Michigan accounts for 49.9% in "electronic stability control."

These steps illustrate how, once we have identified a list of new technologies and their associated phrases, we can build a rich panel dataset of these technologies by identifying their mentions in other text sources. We expand on this theme next.

#### d. Step 4: Cross-reference with job postings

We finally cross-reference our list of technologies with the full text of online job postings, which we source from Burning Glass (BG). BG aggregates online job postings from online job boards (such as indeed.com), employer websites, and other sources into a de-duplicated database.

We employ two datasets from Burning Glass. The first is a standardized dataset (used recently by Hershbein and Kahn, 2018; Deming and Noray, 2020; and Atalay et al., 2020), where each de-duplicated job posting is geo-coded and assigned to a Standard Occupational Classification (SOC) code (the U.S. Bureau of Labor Statistics' system of classifying occupations) and a North American Industry Classification (NAICS) code.<sup>18</sup> The second dataset has thus far received less attention by researchers. It contains the raw unprocessed text of the job postings, which we use to

---

<sup>18</sup> We make extensive use of the former, which are available for 80% of all postings. Industry classifications are available for a more limited 41% of postings. We use industry data only in Section 6. The strings with firm names are available for 66% of all postings.

identify jobs involved with the research, development, production, or use of our technologies. Appendix Figure 3 displays some representative pages from a full BG database entry.

We have data from BG for all available years, 2007 and 2010-2019, a total of roughly 200 million job postings. We drop 2007 jobs from our baseline analysis because Burning Glasses is missing data for 2008-09, though including the 2007 data has little impact on our results.

We associate each posting with a skill level, location, industry, and firm as follows (for details, see Section 1.5 of the Data Appendix): *Skill level*. We construct a skill level for each six-digit SOC code (the most detailed level) given in BG by measuring the share of persons with a college degree, the share of persons with a PhD or a master’s degree, the average wage, and the average years of schooling in the American Communities Survey (ACS 2015 release), using respondents reporting their occupation as in that six-digit SOC code.<sup>19</sup> *Location*. We use the county names provided by BG to uniquely assign job postings to one of the 917 CBSAs in the United States. *Industry*: We allocate a job posting to an industry using the four-digit NAICS code provided by BG.<sup>20</sup> *Firm*: To allocate job postings to firms, we extend the methodology of Autor et al. (2020) and cluster employer strings associated with job postings together on the basis of top search results on Bing.com. For more details on the firm mapping, please refer to Section 3 of the Data Appendix.

To identify job postings associated with each technology bigram, we simply check whether the job posting mentions that bigram and create an indicator variable that is equal to one if it does:

$$Technology\ Job_{i,\tau,t} = 1\{b_\tau \in D_{i,\tau}\}, \quad (1)$$

where  $b_\tau$  is a given technology bigram  $\tau$  associated with one of our new technologies and  $D_{i,\tau}$  is the set of bigrams contained in job announcement  $i$  posted in year  $t$ . In our main specification, we exclude the first and last 50 words of the job posting from this set to avoid picking up mentions of the technology in a description of the firm, as opposed to the task to be performed by the employee, as we discuss below.

To interpret what it means for a job posting to mention a technology, we conduct a human audit of 1,000 randomly selected technology job postings (see Appendix Table 6 for details). As

---

<sup>19</sup> For SOC codes in job postings where we do not find any persons surveyed in the ACS, we match them to the closest available SOC code in the ACS. For example, data for SOC Code 38-1967 were not available, so we match these observations to 38-1960. In total, the dataset includes 837 SOC codes.

<sup>20</sup> NAICS codes typically have six nested levels; the four-digit level is referred to as “industry group.”

expected, the vast majority of mentions relate to a task to be performed by the employee (91% if we trim the first and last 50 words, 80% otherwise). That is, job postings usually mention technologies when the job involves using, producing, or otherwise interacting with the technology. For example, a job ad with mention of “touch screen” (see Appendix Table 6) requires the worker to bend frequently to use the keys on a touch screen to enter text or data in a computer. The remaining mentions are either unspecific (4% in our human audit), for example, mentioning that these technologies are available in the workspace, or refer to the company but not the job (the remaining 5% if we trim the first and last 50 words, 16% otherwise).

For each of our 1,286 technologies, we thus have its year of emergence, a list of pioneer locations where the technology was invented, and a highly granular dataset of job announcements (indexed with a location, industry, occupation, skill level, firm, and year) that involve using, producing, or otherwise interacting with the technology. Most of our analysis focuses on aggregations of these granular data to the technology-time and the technology-location-time levels. Appendix Table 7 provides summary statistics for each level of aggregation. However, the data also open the door for much more granular analyses of job postings for specific firms, locations, and occupations, as we discuss below (see Appendix Table 5 for an example).

Of course, each of the four steps of our data construction can be implemented in different ways, which we highlight when exploring robustness in Section 7. For example, we may choose different thresholds for a technology’s emergence year, include or exclude unigrams or trigrams, and employ various human audits of the technologies identified by our algorithms. While each of these variations result in a slightly different sets of technologies and bigrams, we find they have little effect on our main findings below. It should be noted that a number of studies have used employment data from other sources that we do not explore here to understand the diffusion of technology. Among the most important of these are Tambe and Hitt (2012), Tambe (2014), and Tambe et al. (2020), who measure the skills of U.S. IT workers using resumes from Linked In.

#### e. Disruptive technologies

Figure 2 shows a binned scatterplot of the number of mentions in earnings calls over the number of job postings mentioning each of our 1,286 technologies. It shows two important patterns: First, both variables are highly correlated – the same new technologies that occupy the discussions of managers and investors in earnings calls are also most frequently mentioned in job postings. The



$R^2$  of a fitted regression line is 57.0%. Second, both distributions are heavy tailed (note the logarithmic scale on both axes), so that a small number of technologies drives the vast majority of the mentions in both job postings and earnings calls. The 276 technologies that are mentioned in more than 100 earnings calls account for about 39 million job postings (or about 77 percent of all job postings mentioning any new technology).

For ease of reference, we refer to this highly impactful group as “disruptive” new technologies (Christensen, 1997; Rogers, 1962), in the sense that these new technologies appear to pose enough of an opportunity or challenge to listed firms to take a significant amount of airtime in their quarterly earnings calls. We pay special attention to this group throughout our analysis.<sup>21</sup> It includes all the examples from Table 1 (smartphone, machine learning, hybrid vehicles...). On average, each disruptive technology is mentioned in 141,633 job postings.

The figure also shows examples of “intermediate” and “non-disruptive” new technologies – those with between 10 and 100 and fewer than 10 earnings call mentions, respectively. The former (intermediate) group includes video compression, the pulse oximeter, and the liquid chromatograph. On average, each technology in this group is associated with 26,128 job postings. The latter (non-disruptive) group includes the ultrasonic horn, suction filtration, and NMOS transistors, with on average 1,165 job postings.<sup>22,23</sup>

### 3. Spatial Concentration of New and Disruptive Technologies

We first describe the spatial distribution of innovative activity associated with our new technologies. Table 3 examines the regional concentration of patents that mention new technologies. Panel A lists the top five CBSAs (in column 3) and their share of innovative activity

---

<sup>21</sup> In this sense, our notion of disruption is notably broader than Schumpeter’s (1942) theory of creative destruction, which also requires the devaluation or destruction of existing wealth.

<sup>22</sup> Note that our notion of technology as an “applied technique, product, or process” naturally recognizes the NMOS transistor as a separate technology from the smartphone, even though the latter might contain or even require the former. Similarly, in the context of job postings, there is a clear distinction between a job task requiring use of a smartphone and a job task involving NMOS transistors. In this sense, we are using language, which naturally generates different terms for different technologies that workers and firms interact with, to measure a technology’s economic importance in job postings and earnings calls. These notions of economic importance are thus also quite distinct from broader notions of scientific importance, where understanding electricity and transistors are prerequisites to building smartphones.

<sup>23</sup> Appendix Figure 4 shows the average number of job postings for each category of technology.

(in column 2), broken down by disruptive, intermediate, and non-disruptive new technologies and then for patents mentioning any new technology. It shows two major stylized facts.

First, relative to the distribution of the population and the educated workforce, the development of new technologies is regionally concentrated. Of the 917 CBSAs, the top five collectively account for 33.3% of patents mentioning a new technology. As such, the development of new technologies is notably more concentrated than the distribution of college graduates (22.5%) and the overall workforce (18.9%).<sup>24</sup>

Second, this concentration increases when we condition on patents that mention a disruptive new technology, where the top five CBSAs (San Jose, San Francisco, New York, Seattle, and Boston) represent a striking 42.1% of all patents. By way of comparison, intermediate patents, which less frequently appear in earnings calls, have 34.0% in the top five CBSAs, while non-disruptive patents have 30.4%. *The most commercially impactful innovations thus also have the most geographically concentrated origins.* Notably, only one of the largest CBSAs for disruptive patents – New York – is on the top five list for employment, highlighting how population size is not the primary correlate of disruptive patenting share.

In the same vein, Figure 3 shows the distribution of pioneer CBSAs – the urban areas that account for a majority of *early* patenting of each of our disruptive technologies. Panel A of Figure 3 presents these patterns in map form; and Panel B presents them in a bar chart showing CBSAs' share of all pioneer locations. In Panel B, we combine San Jose-Sunnyvale-Santa Clara, CA and San Francisco-Oakland-Hayward, CA as Silicon Valley, which accounts for 28.7% of all pioneer locations. Jointly, all California CBSAs account for about 41.0% of all pioneer locations. Major cities in the Northeast Corridor, New York-Newark-Jersey City, Boston-Cambridge-Newton, Washington-Arlington-Alexandria, and Philadelphia-Camden-Wilmington, jointly account for 27.6% of all pioneer locations. The top two clusters alone – Silicon Valley and the Northeast Corridor – thus account for 56.3% of all pioneer locations. This result highlights the concentration of disruptive innovative activity within America over the last decades.

These pioneer locations for disruptive new technologies tend to have highly educated workforces and a high density of university activity. For each CBSA-disruptive technology pair (e.g., “smart

---

<sup>24</sup> The table also reveals disruptive technologies are more patent intensive. We can compute 2,434 (1,044,351/429) patents per bigram for disruptive innovations, as opposed to 612 for intermediate and 319 for non-disruptive ones.

phone” and the San Jose CBSA), Appendix Figure 5 presents binned scatter plots of patents mentioning each technology in the ten years prior to the emergence date (per capita, normalized by total CBSA population) and regional characteristics. In the top left panel, the x-axis is the share of workforce with a college degree; in the top right, the share with a post-graduate degree; and on the bottom, the log of university assets per capita. In all cases, there is a strong association between measures of education/university presence and disruptive patents per capita. For example, CBSAs at the 75<sup>th</sup> percentile of university assets per capita file 153% more disruptive patents per capita than those at the 25<sup>th</sup> percentile. We also compute these for non-disruptive technologies. Interestingly, these associations tend to be significantly stronger for disruptive than non-disruptive patenting (although they are positive for both categories). Regions with a greater research university presence or a more educated workforce are thus significantly more likely to be involved in the early development of disruptive technologies.<sup>25</sup>

We show evidence below that this concentration of disruptive innovation in a handful of urban centers engenders large and persistent regional disparities in economic opportunity, as measured by job postings in local labor markets. In this sense, a handful of U.S. locations appear to have a comparative advantage in developing technologies that disrupt firms and labor markets.

#### 4. Diffusion across Regions – Region Broadening and Pioneer Advantage

We next seek to understand the diffusion of new technologies in job postings across regions.

To understand the geographic spread of technology job postings, we define the normalized share of job postings in CBSA  $c$  mentioning a technology bigram  $\tau$  in year  $t$ :

$$Normalized\ share_{c,\tau,t} = \frac{\sum_{i \in c} Technology\ Job_{i,\tau,t} / \sum_i Technology\ Job_{i,\tau,t}}{\#Jobs_{c,t} / \#Jobs_t}. \quad (2)$$

The numerator measures the share of all jobs relating to a given technology  $\tau$  at a given point in time  $t$  that are located in  $c$ ; and the denominator is the share of location  $c$  in the overall U.S. labor market at  $t$ . *Normalized share* <sub>$c,\tau,t$</sub> , therefore, measures the regional over- or under-representation

---

<sup>25</sup> This finding matches the large literature on the geographical concentration of innovation and its connections to university activity, such as Jaffe (1989), Jaffe et al. (1993), Zucker et al. (1998), and Furman and MacGarvie (2007). Moretti (2021) illustrates these effects by examining inventor moves to larger innovation clusters, showing that they experience significant increases in inventive productivity. (This result was hinted at in Forman et al. (2016) as well.)

of job postings associated with each technology bigram relative to the distribution of overall open jobs. Values above one denote over-representation and below one under-representation.<sup>26</sup>

In Figure 4, we present a series of maps displaying the spread of job postings mentioning new disruptive technologies. The blue circles identify the same pioneer locations as in Figure 3, but now superimpose purple dots that show the intensity of the normalized share of job postings relating to these new disruptive technologies 0-5, 6-10, 11-20, and 21-30 years after the disruptive technology's year of emergence. Darker dots correspond to a higher normalized share of jobs.

Two patterns stand out. First, as time goes by, jobs relating to disruptive technologies gradually spread across space (region broadening). Second, there is a remarkable alignment between the CBSAs that pioneer early development in disruptive technologies and the CBSAs that host their early employment. Even after accounting for differences in the size of the local labor market, early employment is strongly concentrated in the same places where the technology was originally developed (pioneer advantage). We next substantiate these two patterns formally.

a. Region broadening

We first examine the overall geographic dispersion of technology job postings. To this end, we calculate the coefficient of variation of the normalized share of technology job postings by dividing the standard deviation of *Normalized share*<sub>*c,τ,t*</sub> across locations *c* in year *t* by its mean in year *t* for each technology bigram  $\tau$ .<sup>27</sup> If technologies are uniformly spread out across CBSAs, then the normalized share takes a value of 1 for each CBSA, and the coefficient of variation calculated across CBSAs is 0. The average coefficient of variation in our sample of disruptive technology bigrams is 5.70, which suggests that technology job postings relating to disruptive technologies are on average highly concentrated compared to, for example, the coefficient of variation for the normalized share of the local population that holds a college degree (2.90).

Using a regression framework, Table 4 examines the evolution of this coefficient of variation over the technology's life cycle. Panel A of this table reports results from regressions of the form:

$$CV_{\tau,t} = \alpha_0 + \beta_{RB}(t - t_{0,\tau}) + \delta_{\tau} + \varepsilon_{\tau,t}, \quad (3)$$

<sup>26</sup> Throughout, we cap this variable at the 99<sup>th</sup> percentile of non-zero observations.

<sup>27</sup> Appendix Table 7 summarizes the data used in this and subsequent regression analyses.

where  $CV_{\tau,t}$  is the coefficient of variation across CBSAs for technology bigram  $\tau$  in year  $t$ , and  $(t - t_{0,\tau})$  is the number years since emergence of technology bigram  $\tau$  in year  $t_{0,\tau}$  (capped at 30 years, given we have little data for technologies older than 30 years).  $\delta_\tau$  denotes a full set of technology bigram fixed effects, which we constrain to sum to zero, so that the intercept  $\alpha_0$  measures the average coefficient of variation at the year of emergence.<sup>28</sup> The slope coefficient,  $\beta_{RB}$ , measures the speed of decay of this concentration with each passing year since emergence. Panel A, column 1 reports estimates for disruptive technologies. Columns 2 and 3 report results for all new technologies, without and with bigram fixed effects, respectively. Throughout, we cluster standard errors at the technology level (the unit of observation is a technology bigram).

In column 1, we find that the coefficient of variation, on average, decreases by 0.068 (s.e.=0.026) for every year since its emergence year. Job postings mentioning disruptive technologies slowly, but significantly, spread across space. On average, at the year of emergence of a disruptive technology, the coefficient of variation is 5.58 (significantly greater than 0). With each additional year since emergence, on average, this coefficient of variation decreases by 0.068 points (or 1.22%). Taking these estimates at face value suggests that disruptive technology job postings take 82 years to fully disperse across the United States.

Figure 5 shows this pattern graphically using a binned scatterplot: a technology’s job postings are geographically highly concentrated in the early years after its emergence. Within 30 years, this geographic concentration drops by about a third (36.6%). Interestingly, the figure also shows this process of spread, measured in the pooled set of technologies, is close to linear in the data.<sup>29</sup>

In column 2 of Table 4, we show this pattern is almost identical when we include the long tail of intermediate and non-disruptive technologies in our sample. Appendix Table 18 tests explicitly for differences in the rate of spread between disruptive and non-disruptive technologies, finding few

---

<sup>28</sup> Because the coefficient of variation, as well as some of the other constructed moments used in the following tables, become noisy with insufficient data, we take steps in the regressions to down-weight technologies that are mentioned in relatively few job postings. First, we weight observations by the square root of the total number of job postings mentioning that technology, capped at 100, meaning that technologies with more than 10,000 postings receive full weight, while those with less than 10,000 postings are weighted by their square root. Second, in regressions not limited to disruptive technology bigrams, we exclude technology bigrams with less than 1,000 job postings. In practice, these adjustments have little impact on our estimates (see Section 7).

<sup>29</sup> Appendix Table 8 includes a quadratic term and shows it is indistinguishable from zero. Additional specification tests suggest the relationship is closer to linear than log in the data. Consistent with the literature on S-curves, which studies the speed of adoption of a given technology (Griliches, 1957), we do find significant concavity in the rate of spread when conditioning only on variation *within* technologies (column 3 of Appendix Table 8).

significant differences. In column 3 of Table 4, when conditioning only on within-technology variation, we find a somewhat faster rate of spread. With each additional year, the coefficient of variation falls by -0.153 (s.e.=0.012) or 1.85% – implying 54 years to full dispersion.

Panel B of Table 4 shows similar results (following the same specification as column 3 of Panel A) using alternative measures of geographic concentration as dependent variables: the mean normalized share of a technology’s job postings in the top five CBSAs relative to the mean normalized share across all CBSAs, the percentage of CBSAs with a normalized share of a technology’s job postings of less than 10% (that is, the representation of CBSAs with almost no activity associated with that bigram), and the sum of squared deviations of the normalized share from 1 (similar to the Herfindahl-Hirschman Index). The consistent pattern is for a slow decline of concentration, however measured: all three measures fall with time, but again imply time periods in excess of 50 years to full dispersion. This is a strikingly slow rate of convergence, given that the typical political cycle is around five years, and most Americans work for less than 50 years.

b. Pioneer advantage

Table 5 formally explores the second pattern: pioneer advantage. We quantify the advantage that pioneering regions (CBSAs that account for a majority of the initial patenting in a technology) retain in that technology’s job postings, even as region broadening occurs. Panel A reports results from the specification:

$$Normalized\ share_{c,\tau,t} = \alpha_0 + \beta_P Pioneer_{c,\tau} + \beta_D Pioneer_{c,\tau}(t - t_{0,\tau}) + \delta_c + \delta_\tau + \delta_t + \varepsilon_{c,\tau,t} \quad (4)$$

where  $Pioneer_{c,\tau}$  is a dummy variable denoting the pioneer status of the CBSA;  $\delta_c, \delta_\tau, \delta_t$  denote CBSA, technology bigram, and year fixed effects respectively. Columns 1 and 2 examine job postings mentioning disruptive technologies; columns 3 and 4 show results for all technologies.

In column 1, we see that the pioneer locations associated with a disruptive technology retain a significant pioneer advantage on average. More specifically, the normalized share of technology job postings is 31.1 percentage points higher in its pioneer locations on average throughout the sample period. Column 2 shows that this advantage is much larger in the year of emergence (108.4 percentage points), but then decreases significantly over time, on average by 3.2 percentage points per year or 2.9% (0.032/1.084). The initial advantage of the pioneering locations for job postings

relating to the disruptive technologies they develop thus lasts for decades, with an implied 34 years to zero advantage.<sup>30</sup>

In column 3, we include all technologies and find an almost identical pattern – albeit with a somewhat larger point estimate for the pioneer advantage at the year of emergence of 1.321 (s.e.=0.254). In column 4, we look at technology job postings in the neighborhood of pioneer locations by adding a dummy for CBSAs within 100 miles of a pioneer location, *Pioneer Neighbor<sub>c,t</sub>*, and its interaction with the number of years since the emergence of the technology. The estimates suggest that some of the pioneer advantage spills over to these adjacent communities, with a 15.8 percentage point higher normalized share in the year of emergence. Again, this advantage appears to decay over time, though the decay is not statistically distinguishable from zero.

### c. Mechanisms

Given the extreme regional concentration of disruptive technologies' pioneer locations, and the long-term advantage in jobs these regions appear to enjoy, a key question is *why* this advantage appears to be so persistent. We take two steps to better understand the mechanisms behind this persistence: First, we examine the skill requirements of the jobs spreading across space. Second, we analyze the words around those in which the new technology is mentioned in the job posting to learn about whether the job is involved with developing or using the new technology.

#### *Pioneer advantage in high- vs low-skill jobs*

We first analyze differential rates of spread of high- versus low-skill jobs relating to disruptive new technologies. To compute a job posting's skill requirement, we use the 6-digit SOC code allocated to the job posting by Burning Glass and assign it the average level of college education respondents report in the 2015 ACS for that occupation.<sup>31, 32</sup>

Columns 1 and 2 of Panel A of Table 6 report results from the specification:

---

<sup>30</sup> In Appendix Table 9, we test the robustness of our results to the addition of interacted fixed effects. We find that decay rates of pioneer advantage are similar across these specifications.

<sup>31</sup> As an example, Appendix Table 10 shows the list of top occupations by share of job postings for some of our top technologies. Also see Section 2 of the Data Appendix for details.

<sup>32</sup> The BG data also includes an indicator for college requirement for a subset of observations. However, since this subset is quite limited, we prefer using SOC codes to generate this variable.

$$\log(CV_{\tau,t}^s) = \alpha_0^s + \gamma_1^s(t - t_{0,\tau}) + \delta_\tau + \varepsilon_{\tau,t}, s \in \{H, L\} \quad (5)$$

where  $CV_{\tau,t}^s$  is the coefficient of variation of the normalized share of technology job postings across CBSAs, as in Section 4.a, calculated separately for  $s \in \{H, L\}$  – high-skill jobs ( $H$ ) and low-skill jobs ( $L$ ). For the purposes of this exercise, we define high-skilled jobs as those which are classified in occupations with more than a 60% college-educated share in the 2015 ACS (28.4% of all jobs on BG) and low-skilled jobs as those with under a 30% share (42.5% of all jobs). Column 1 reports results for high-skill jobs, and column 2 reports results for low-skill jobs. To facilitate the direct comparison of differential rates of spread between these two types of jobs, we take logs of the dependent variable so that the slope coefficient is now directly informative about the percentage decline in the coefficient of variation per year.<sup>33</sup>

We find that the geographic concentration of low-skill jobs (in column 2) decreases 1.2 percentage points or 40% ( $0.042/0.030 - 1$ ) faster than that of high-skill jobs (in column 1). This difference is statistically significant at the 1% level. Figure 6, Panel A shows this differential decay graphically, this time also including across-technology variation (without technology bigram fixed effects). Again, low-skill technology job postings spread at a significantly faster rate.

This pattern is similarly prominent when analyzing pioneer location advantage. In columns 1 and 2 of Panel B of Table 6, we repeat the regression specification in column 2 of Table 5, but now separate between high- and low-skill jobs (all definitions are as above). We find that the pioneer advantage in a technology's job postings is significantly more persistent for high-skill jobs than for low-skill jobs. While the former decays at 2.4% per year, the latter erodes at a relatively fast 3.6%. These estimates imply it takes 41 years for a pioneer location's advantage in high-skill jobs to erode, whereas that for low-skill jobs lasts only 28 years. (Both results are again almost identical when we repeat these estimations for all of the new technologies in Appendix Table 19.)

Taken together, this evidence suggests that the overall geographic spread of technology jobs is driven by low-skill jobs, while high-skill jobs take significantly longer to spread across space. That is, the pioneer locations involved with the early development of a technology tend to retain a significant and very long-lasting advantage in high-skill job postings relating to that technology.

---

<sup>33</sup> Results are almost identical when we return to the linear specification of Table 4 and compare the ratios of the slope coefficients to constant terms. See Appendix Table 11.



### *Research, Development, and Production*

While there are a number of hypotheses that can be offered for these patterns, our text-based methodology allows us to look carefully at one leading explanation: the movement of new jobs from technology research, development, and production (RDP) to technology use.

The text in the job announcements contains rich information to distinguish these two types of jobs. For example, a job posting involved with a technology's RDP might state "*you will be designing the graphics module for our **virtual reality** training system,*" while one involved with a technology's use might read "*the role will involve assisting customers and selling tickets from your **smart tablet** in the entrance of the cinema.*" (Additional examples in Appendix Figure 6.)

To systematically identify the cases that involve RDP of new technologies, we use an iterative procedure that combines an unsupervised learning algorithm with some human judgment to identify word patterns associated with RDP job postings. The first step is developing a set of plausible keywords (generated by the authors) that are commonly used when describing positions relating to the RDP of new technologies ("research and," "and develop," "and development," "customization of," "to build," and "to design"). We then use an embedding vector algorithm trained on earnings calls to identify other phrases (unigrams and bigrams) that were used in similar context to these keywords – in effect, using the embedding model like a custom-trained thesaurus.<sup>34</sup> For each of these suggested phrases, we examine ten excerpts from job postings to check for false positives. We then add to our initial list those suggested phrases that had at least eight true positives (no more than two false positives). After updating the list, we go through the steps again iteratively – now asking the embedding model for phrases proximate in vector space to the union of already selected phrases – until we have exhausted all useful suggestions that meet the threshold of eight out of 10 true positives. Appendix Table 12 lists the full set of selected phrases.

Using this classification, we systematically flag all job postings that mention a new technology within 15 words of one of our RDP keywords and categorize all others under "use." To verify the

---

<sup>34</sup> Specifically, we use the Word2Vec Python package Gensim trained on earnings calls (sourced as noted above) from 2002 to 2019. For the training process, we used the default parameters: 200 dimensions, ignoring words that appear fewer than 50 times, and a context window of 15 words. We train on earnings calls, instead of job postings, because this type of language model tends to perform poorly when trained on short texts.

accuracy of the resulting classification, we conduct a human audit of 1,000 randomly sampled technology job postings. We assign team members to read and classify these job postings into either RDP or use of the associated technology. In this random sample, we are able to correctly classify 63.1% of technology RDP postings and 68.1% of technology use postings. With this distinction in hand, we calculate the coefficient of variation of the normalized share of technology job postings for each disruptive technology and year separately for the RDP and use job postings. Columns 3 and 4 of Table 6 (Panel A) examine the differential spread of these two different types of technology job postings, estimating region broadening separately for each group. Again, we see large differences: technology-using job postings spread out 94.4% ( $=0.035/0.018-1$ ) faster than postings that involve technology research, production, and development jobs. This difference is significant at the 1% level.

We find a similar pattern for the pioneer advantage in RDP job postings. Columns 3 and 4 of Panel B in Table 6 re-estimate regression specification (4) and calculate the advantage of pioneer CBSAs in technology job postings that involve the RDP and use of new technologies. We find that pioneer advantage in job postings involving the use of new technologies is smaller initially (with a constant term suggesting 147.2% more such jobs in the pioneer location in the year of emergence) and dissipates significantly over time ( $-0.038$ ,  $s.e.=0.014$ ). By contrast, RDP job postings are more concentrated in pioneer locations initially (213.8% higher in the year of emergence), and the decay rate is statistically indistinguishable from zero (though negative and in a similar range as other estimates in the table). (see also Figure 6, Panel B.)

Taken together, these findings suggest that technologies remain highly concentrated in their research, development, and production in the original pioneer location, using highly skilled employees for these activities, but spread out in their application, where lower-skilled employees are utilized. To consider the example of smart phones, these continue to be developed primarily in Silicon Valley by graduate school- and PhD-educated employees, but jobs involving their use have spread out across the U.S., including positions for sales, repair, maintenance, and utilization, often undertaken by non-college-educated employees. That is, pioneer locations that initially developed a technology appear to retain a long-term advantage in high-skilled jobs, because activities relating to the technology's RDP remain in that location for long periods of time.

## 5. Skill Broadening

We next turn to examining the skill bias of technology job postings over time. We find a significant high-skill bias in new technologies initially. Over time, the share of lower-skilled job postings mentioning the technology increases, albeit at a relatively slow rate.

We compute the average skill requirement of job postings associated with a particular technology bigram at a point in time by examining the occupational composition of these job postings:

$$Skill_{\tau,t} = \frac{\sum_o N_{o,t}^{\tau} \chi_{o,2015}}{\sum_o N_{o,t}^{\tau}} \quad (6)$$

where  $N_{o,t}^{\tau}$  is the number of Burning Glass job postings mentioning bigram  $\tau$  that are in SOC code  $o$  at time  $t$ , and  $\chi_{o,2015}$  is the average skill level for occupation  $o$ , as measured by the 2015 ACS. For example, if for a technology bigram  $\tau$  in year  $t$  all associated job postings are in an occupation  $o$ , then its skill level equal to the average skill level of workers in occupation  $o$  in the ACS.

Table 7 uses a regression framework to describe the evolution of the skill level of job postings associated with new technologies. The specification is identical to equation (3):

$$Skill_{\tau,t} = \alpha_{0,SB} + \beta_{SB}(t - t_{0,\tau}) + \delta_{\tau} + \varepsilon_{\tau,t}, \quad (7)$$

but now we use the average skill required for jobs associated with technology bigram  $\tau$  in year  $t$  as the dependent variable. The intercept  $\alpha_{0,SB}$  denotes the average skill level of the technology's job postings in its year of emergence,  $t_{0,\tau}$ . The slope ( $\beta_{SB}$ ) denotes this skill level's average speed of decay with each passing year since emergence. Column 1 of Panel A reports results for disruptive new technologies. Columns 2-4 again include all new technologies.

In column 1, we find that, on average, 57.1% of job postings mentioning a new technology require a college degree at the year of emergence of the technology. As such, jobs associated with a new technology are significantly skill biased, particularly when compared with the share of the U.S. workforce that holds a college degree – about one third. At the same time, this skill content of a technology's job postings is significantly downward sloping over time. With each additional year since emergence, it falls by 0.228 (s.e.=0.092) percentage points on average, implying a rate of skill broadening of 0.40% ( $=-0.228/57.078$ ) per year.

Figure 7 shows this evolution graphically using a binned scatterplot. Although the pattern of skill broadening is clearly visible, it is worth noting that 30 years after the year of emergence, the

average college requirement is still 50.2%, far above the average rate of college attainment in the U.S. population, as noted above. In this sense, new technologies persistently generate a disproportionate share of employment opportunities for high-skill workers for very long periods of time. Column 2 of Table 7 shows almost identical results when we include all technologies in the sample.

One possible concern with these results is that the types of jobs advertised online (as opposed to in printed newspapers) could be changing over time.<sup>35</sup> To address this concern, Column 3 shows the coefficient of interest is almost unchanged when including time fixed effects (-0.218, s.e.=0.100), so that our findings cannot be explained by an increasing share of low-skilled jobs being advertised online. Appendix Table 13 expands on this theme by estimating skill bias and broadening separately for two sub-samples (2010-2015 and 2016-2019), with almost identical results in each case.

In column 4, we introduce technology bigram fixed effects and now find a larger negative slope (-0.493, s.e.=0.036), but also a larger constant term (63.898, s.e.=0.840).<sup>36</sup> Taken at face value, the two estimates imply that new technologies take 68.08 years to reach the average level of college education among the U.S. workforce (30.3% in the 2015 ACS). In other words, the skill bias of a given new technology on average takes several generations to dissipate.

Panel B of Table 7 repeats this estimation using alternative measures of skill. It shows that, in the year of emergence, jobs in a new technology on average require 15.5 years of schooling, 22.6% of them require a post-graduate degree, and they pay an average wage of \$75,521 (measured in 2015 dollars). All three skill indicators again decay significantly over time, at rates that would imply 77.6, 78.0, and 69.7 years to reach the average years of schooling, rate of post-graduate education, and wage of the U.S. population reported in the ACS.

---

<sup>35</sup> Appendix Figure 9 describes the overall volume and the composition of Burning Glass (BG) job postings over time. Panel A shows that BG job postings have increased about one-to-one with job postings captured in the U.S. Bureau of Labor Statistics' Job Openings and Labor Turnover Survey (JOLTS). Panel B shows that the average skill level associated with BG job postings has fallen over time at about 0.7% per year. Panels C and D show that the volume of BG job postings by occupation (pooled across years and by year) is associated one-to-one with employment observed in that occupation, indicating that BG has been consistently representative of U.S. employment.

<sup>36</sup> Note that, due to the linear form of our estimating equation, the within-technology-and-time variation is effectively degenerate, so that we cannot simultaneously introduce technology and time fixed effects. In this sense, there is no way of distinguishing cohort from time effects, as is common in such analyses (see Hall et al., 2007).

All of these variations show (1) that job postings mentioning new technologies are strongly high-skill biased initially and (2) this skill bias decays significantly over time, albeit at a relatively slow rate, so that the skill bias of jobs associated with new technologies persists for multiple decades.

Both findings intersect with important branches of the literature studying the relationship between technology and inequality. First, they show direct evidence of the high-skill bias of new technologies, adding to a large literature that infers this skill bias from observed wage premia (e.g., Katz and Murphy, 1992). The findings suggest in a dramatic way that new technologies contribute to persistent inequalities between high- and low-skilled workers and, because pioneer locations of disruptive technologies are highly concentrated, also engender persistent inequalities across space. In this sense, innovation has a profound effect on regional disparities in economic opportunity.

Second, our finding of skill broadening provides direct evidence for a key assumption in the literature on automation: that the comparative advantage of high-skill workers in a new task erodes as the technology matures, pulling lower-skilled workers into working with a new technology over time. It is this key assumption that leads to balance in the “race between man and machine” in Acemoglu and Restrepo (2018) and the related literature. Our evidence suggests this skill broadening indeed occurs in the data.

#### a. Mechanisms

Given these results, a key question is why skill broadening occurs in practice. The literature has suggested at least two, possibly complementary, channels. The first is standardization of new technologies – where research and customization become less important as new technologies mature and become standardized. That is, the new technology evolves over time into a standardized form that can more readily be used by less educated workers (Acemoglu et al., 2012; Acemoglu and Restrepo, 2018). The second is training – over time, less educated workers may acquire experience or training that allows them to use new technologies, even if they do not have high levels of formal education (Nelson and Phelps, 1966; Galor and Moav, 2000).

Again, analyzing the context of the mention of the new technology within a given job posting can shed some light on these mechanisms. To this end, we use our keyword-based approach to systematically flag those job postings that mention a given new technology in conjunction with a requirement of training or experience with that technology (starting with seed phrases “training in,” “knowledge of,” “experience with,” “familiar with,” “knowhow of,” and “proficiency in”);

and we again use the same iterative procedure combining our embedding vector model with human reading to settle on a list of keywords (shown in Appendix Table 14).

Appendix Figure 7 shows the proportion of RDP jobs declines significantly over time, so that more mature technologies have a lower share of jobs involved with RDP. At the same time, these RDP technology jobs skew heavily on the side of higher college requirements. That is, as one might expect, the most educated workers are the ones conducting research, development, and production of technologies. These plots are thus consistent with theories that, as technologies mature, their jobs move from those involving RDP to ones that simply use the technology.

By contrast, Appendix Figure 8 shows that the share of technology job postings that require training or experience in that technology is relatively stable throughout the life cycle of the technology, without a significant upward or downward trend as the technology ages. At the same time, these training requirements are positively, not negatively, associated with college requirements, so that training in a new technology and formal education appear to be complements, not substitutes in our data.

To assess to what extent these two channels can account for new technologies' skill broadening over time, Table 8 separately adds both as controls, to assess to what extent their inclusion can attenuate the estimated coefficient,  $\beta_{SB}$ . Column 1 reproduces our estimate from column 1 of Panel A, Table 7 for comparison (-0.228, s.e.=0.92). Column 4 shows that controlling for the share of RDP jobs in the same technology attenuates this estimate by about a quarter to -0.177 (s.e.=0.075). Columns 2 and 3 shows similar, albeit somewhat smaller, attenuations when controlling separately for the share of R&D job postings and the share of job postings relating to production.<sup>37</sup> We conclude that technologies' transition from a focus on RDP towards a focus on use can account for part of the skill broadening we observe in the data.

By contrast, column 5 shows that controlling for the share of that technology's jobs that require training or experience in the technology results in no attenuation whatsoever of our estimate of  $\beta_{SB}$ . In this sense, changes in training and experience in the technology cannot account for the pattern of skill broadening observed in the data.

---

<sup>37</sup> For the sub-topic of research and development we start with the seed phrases "research and," "and develop," "and development," and "customization of" – a subset of our RDP seed keywords above – and proceed in the same manner. The remainder of the RDP keywords constitute the "produce" category.

We tentatively conclude that training and experience does not appear to be a substitute for formal education when it comes to required qualifications for jobs in new technologies, as measured in job postings. Instead, some of the observed skill broadening can indeed be accounted for by standardization of the technology over time.

## 6. Diffusion across Occupations, Industries, and Firms

Before exploring the robustness of our main findings, we highlight the power of the data that we have developed to also characterize the spread of new technologies across other dimensions.

To assess the rate at which technologies spread across occupations, firms, and industries, we extend the definition of *Normalized share* $_{c,\tau,t}$  to NAICS four-digit industries, SOC six-digit occupations, and firms for each technology ( $\tau$ ) and time ( $t$ ), calculating the normalized share of job postings in each industry, occupation, and firm that mention a given new technology.<sup>38</sup> We then measure the coefficient of variation of *Normalized share* $_{c,\tau,t}$  across the segments.

Because the number of firms posting job advertisements online expands over time, we stratify our firm-technology-year sample by including only firms that post at least one job in each of our sample-years, before calculating the coefficient of variation.<sup>39</sup> This step focuses attention on 10,231 larger firms, which on average post 1,628 jobs per year, effectively excluding variation coming from small and medium-sized businesses.

*Spread across firms, occupations, and industries.* Table 9, Panel A shows the results of a regression of the coefficient of variation calculated for each technology ( $\tau$ ) and time ( $t$ ) on the year since emergence. Column 4 shows our already established results for locations for comparison.<sup>40</sup> We find that while there is a decline in concentration as measured by coefficient of variation for all four segments, there is a relatively (and significantly) larger decline across locations and firms (columns 4 and 3) than across industries and occupations (columns 2 and 1).

---

<sup>38</sup> While the former two variables are included in the BG data (in each case, we use the finest level of disaggregation available from BG), the latter relies on our own matching algorithm described in Section 2.

<sup>39</sup> Hershbein and Kahn (2018) discuss this fact in some detail. The general increase in coverage of the BG data over time should not affect any of our main results. We discuss robustness to various weighting schemes in detail in Section 7.

<sup>40</sup> In order to avoid calculating coefficients of variation for unreasonably sparse data, we only keep those technology x year observations with at least 100 postings with industry coverage. This issue arises because BG provides NAICS codes for only 41% of all postings, as noted above.

While the coefficient of variation declines on average by 1.8% and 1.6% per year for CBSAs and firms, respectively, the corresponding declines are 0.7% and 0.4% for occupations and industries, respectively.<sup>41</sup> In fact, in Column 1, this rate of decline across industries is statistically indistinguishable from zero.

*Advantages for pioneer firms and industries.* Following our procedure for pioneer locations, we define pioneer industries and firms for each technology as those with the most assigned patents in the ten years after the technology’s emergence year that collectively account for 50% of the matched patents in a given new technology. (See Section 3 of the Data Appendix for details on how we match patents to large firms and industries – matching patents to occupations makes little sense, so that we do not calculate pioneer occupations.) Appendix Table 15 shows the top pioneer industry for a selection of technologies. For example, the top pioneer industry for “hybrid electric vehicle” is “motor vehicle manufacturing” (accounting for 56% of early patents).

In Table 9, Panel B, we explore the initial hiring advantage of pioneer firms and industries by estimating specification (4) for these additional dimensions. The table shows that pioneering firms have a strong initial advantage in job postings, with a 2,093% higher normalized share of job postings at the year of emergence for pioneer firms. Over time, this advantage again degrades significantly, at a rate of 2.3% per year. Consistent with the results in Panel A, this rate of decline is statistically indistinguishable from zero for industries.

Taken together, this evidence suggests new technologies initially generate hiring that is highly localized by location, firm, and industry. Over time, this hiring disperses, particularly across locations and across firms. Looking in more depth on a within-firm basis at the dynamics around the location of innovation and job creation is a fertile avenue for future exploration.

## 7. Robustness Checks and Extensions

Finally, we conduct a broad range of robustness exercises to assess to what extent judgments we have made could have affected our primary results: “concentration in the development of

---

<sup>41</sup> Decay rate across CBSAs is 0.017 (0.005), 0.012 (0.003), and 0.004 (0.002) higher than industries, occupations and firms, respectively. These coefficients are statistically significant at the 1%, 1%, and 10% level, respectively.



disruptive technologies,” “region broadening,” “pioneer-location advantage,” “skill broadening,” and “differential region-broadening by skill level.”

To this end, we first re-trace our four steps of data construction to reexamine each of the main decisions we made in this automated process. In each case, we alter one aspect of the process, re-create our entire dataset, and re-run our main analyses. Table 10 reports the main estimates of interest, where the first line of each panel reproduces the results of our baseline specification for comparison.

*Influential patents (Step 1 in Section 2).* When isolating new bigrams associated with influential innovations, we retained only those that appear in patents accumulating a total of at least 1,000 weighted citations. Having some such threshold is necessary to maintain computational feasibility (to avoid having to cross-reference 1.5 million novel bigrams to Wikipedia and our other text sources). However, Panel A of Table 10 shows our results are almost invariant to altering this threshold. The panel shows four variations, with cutoffs ranging from 1,250 to 2,000, each producing almost identical results.<sup>42</sup>

*Phrase length (Step 1 in Section 2).* Our methodology easily extends to including trigrams, in addition to bigrams in the analysis. Repeating our steps 1-4 for trigrams adds 328 technology trigrams. 262 of these simply add another phrase to the set of bigrams already associated with a given technology (Wikipedia title) in our data. Perhaps the only substantive additions are “real time communications” and “injection molding machine” (see Appendix Table 16).

Adding unigrams is slightly more complicated due to their sheer number (about 2 million pass the threshold of 1,000 cite-weighted patents, simply because unigrams are more frequent than bigrams). To keep the number of candidate unigrams manageable, we focus on those with more than 100 mentions in earnings calls. Doing so adds 200 new technology unigrams, 53 of which again simply add another phrase to the set of phrases associated with a given disruptive technology already identified in our bigram-based analysis. Appendix Table 17 shows examples among the 147 remaining unigrams. Overall, as expected, the unigram-based approach appears significantly noisier, with some clear false positives (“billable,” “internets”) and names in the mix (“USPS”).

---

<sup>42</sup> The reason for this stability is apparent in Appendix Figure 10, which shows a strong correlation between the number of cite-weighted patents and job postings in which a technology is mentioned across all novel bigrams (i.e., including bigrams with few cite-weighted patents). That is, variations in our minimum citations cutoff will on average tend to remove technologies that have little traction in the labor market.

Nevertheless, broadening our approach in this way also yields some substantive additions, including, for example, “mRNA” and “Bluetooth.”

Re-running our analyses including these sets of unigrams and trigrams again has no significant effects on our results.

*Human audit (Step 2 in Section 2).* Rather than relying fully on our Wikipedia filter to determine whether or not a novel and influential bigram describes a technology (as opposed to increasingly visible problems or management techniques), we also conducted a human audit, where team members read through each Wikipedia title – technology bigram pair and removed all of those where the match appeared erroneous (e.g. “OS-level virtualization” matched to “programs running”) and those where either the Wikipedia title or the technology bigram did not describe a technology according to the team member’s judgment (e.g. “adverse event”). Appendix Table 3 marks each of the disruptive technologies dropped under this audit (altogether 63 of 276 disruptive technologies). Doing so again has a negligible effect on our estimates.

*Emergence years (Step 3 in Section 2).* Our baseline approach to defining a technology’s emergence year requires that technologies are mentioned in at least 100 cite-weighted patents prior to their year of emergence. Two variations in Panel D loosen (one cite-weighted patent) and tighten (200 cite-weighted patents) this requirement. A third variation abandons this approach altogether and instead fixes the emergence year as the first year in which the technology reaches 50% of its maximum cite-weighted patents achieved by a technology bigram in our sample. All of these variations again have a negligible effect on our main results.

Note that each of these variations in the robustness checks above alters the list of new technologies we uncover in small ways. For example, “fracking” may only show up in our data if we explicitly allow for unigrams, in addition to bigrams. Similarly, requiring 2,000 rather than 1,000 cite-weighted patents before including a new bigram from patents in our first step of data construction will obviously shorten the list of new technologies we produce. Our measure of success is thus not to always produce the one true list of new technologies that arose in the past 40 years. Such an absolutely true list does not exist. Instead, the key is that our language-based approach produces a list of technologies that is representative of new technologies in a statistical sense. The fact that all of the variations above produce very similar econometric results is evidence that we meet this bar.

*Alternative weighting schemes (Step 4 in Section 2).* Because the coefficient of variation, as well as other constructed moments at the technology-time level, become noisy with insufficient data, our baseline specifications down-weight technologies that are mentioned in relatively few job postings. Panel E repeats all analyses with (i) unweighted regressions, (ii) without the requirement of a minimum number of mentions in job postings, and (iii) with weights proportional to the natural logarithm of the number of job postings associated with the technology. Finally, we re-run our entire analysis after collapsing technology bigrams at the technology (Wikipedia title) level. Again, none of these variations materially affect our results.

*Representativeness of the BG sample.* To further address any concerns relating to the possibly changing composition of the BG data over time, Appendix Table 13 shows additional variations of Table 7, Panel A, column 2 where we (i) include the 2007 data and (ii) estimate our baseline coefficient separately for two sample periods (2010-2015 and 2016-2019).

*Standard errors.* We also explore the robustness of the results relative to the treatment of the standard errors. These examine again the four regressions that were analyzed in Table 10. We explore in Table 11 the impact on the standard errors of different clustering approaches: clustering the observations not by associated Wikipedia entries (“technologies”), but rather by the individual technology bigram, the year, and (in the case of the regression from Table 5) the CBSA, state, and the interaction between the CBSA and the associated Wikipedia entry. We also present bootstrapped standard errors, drawn from 1,000 replications with replacement. The changes have little effect on the significance of the results.

## 8. Conclusion

Policymakers in many parts of the world devote enormous energy to fostering nascent technologies, ranging from efforts to support academic research to luring start-ups from other cities and nations. Such infant industry strategies are often predicated on the notion that early advantages in innovation and employment will yield lasting benefits for regions, particularly in the form of high-quality employment.

Using the full text of patents, job postings, and earnings conference calls, we introduce in this paper an approach to understand which new technologies affect jobs and businesses and to trace

their diffusion across regions, industries, occupations, and firms. We can then map the spread of new technologies in these dimensions, focusing on the hiring associated with each important innovation.

We highlight first that the locations where economically highly impactful technologies are developed are geographically highly concentrated, with a handful of urban areas contributing the bulk of the early patenting and early employment within new technologies. One striking figure is that 56% of the pioneering locations for the most economically impactful technologies are in two parts of the U.S. – Silicon Valley and the Northeast Corridor. Second, despite this initial concentration, jobs relating to new technologies spread out geographically. But this rate of diffusion is extremely slow, happening over several decades rather than in just a few years. Locally developed technologies continue to offer long-lasting benefits for jobs in their pioneer locations for multiple decades. Third, jobs relating to new technologies are highly skill biased – 57% of the initial jobs associated with a given new technology require a college degree. Over time, the mean required skill levels of the new jobs decline, albeit at a very slow pace. Fourth, low-skill jobs associated with the use of a given new technology spread out geographically significantly faster than high-skill ones, so that the pioneer locations where the technology was invented host a disproportionate share of high-skilled jobs relating to that new technology for several decades after its year of emergence.

Combined with the extreme spatial concentration of the economically most impactful innovation, this pioneer advantage engenders large and persistent regional disparities in economic opportunity, giving a handful of U.S. locations a lasting advantage in high-skill jobs.

Beyond these core results of our analysis, the development and spread of new technologies are key objects of interest in multiple fields of economics. As we suggest in Section 6, these techniques developed here should have applications for studies of firm-level technological adoption and implementation. More generally, we hope the text-to-data techniques we develop and data that we provide as part of this paper may prove useful in addressing a range of additional research questions in the study of economic growth, inequality, entrepreneurship, and firm dynamics.

## References

- Abis, Simona, and Laura Veldkamp. "The changing economics of knowledge production." Working paper 3570130, SSRN (2020).
- Acemoglu, Daron, and David Autor. "Skills, tasks and technologies: Implications for employment and earnings." In Orley Ashenfelter and David Card (editors), *Handbook of Labor Economics*. New York, Elsevier, volume 4, chapter 12, pp. 1043-1171 (2011).
- Acemoglu, Daron, Gino Gancia, and Fabrizio Zilibotti. "Competing engines of growth: Innovation and standardization." *Journal of Economic Theory* 147 (2012): 570-601.
- Acemoglu, Daron, and Joshua Linn. "Market size in innovation: Theory and evidence from the pharmaceutical industry." *Quarterly Journal of Economics* 119 (2004): 1049-90.
- Acemoglu, Daron, and Pascual Restrepo. "The race between man and machine: Implications of technology for growth, factor shares, and employment." *American Economic Review* 108 (2018): 1488-1542.
- Acemoglu, Daron, and Pascual Restrepo. "Robots and jobs: Evidence from US labor markets." *Journal of Political Economy* 128 (2020): 2188-2244.
- Agrawal, Ajay, Joshua Gans, and Avi Goldfarb (editors). *The Economics of Artificial Intelligence: An Agenda*. Chicago, University of Chicago Press (2019).
- Akerman, Anders, Ingvil Gaarder, and Magne Mogstad. "The skill complementarity of broadband internet." *Quarterly Journal of Economics* 130 (2015): 1781-1824.
- Arora, Ashish, Sharon Belenzon and Lia Sheer. "Knowledge spillovers and corporate investment in scientific research." *American Economic Review* 111 (2021) 871-898.
- Atalay, Enghin, Phai Phongthientham, Sebastian Sotelo, and Daniel Tannenbaum. "The evolution of work in the United States." *American Economic Journal: Applied Economics* 12 (2020): 1-34.
- Autor, David H., David Dorn, Gordon H. Hanson, Gary Pisano, and Pian Shu. "Foreign competition and domestic innovation: Evidence from US patents." *American Economic Review: Insights* 2 (2020): 357-374.
- Autor, David H., Lawrence F. Katz, and Melissa S. Kearney. "Trends in US wage inequality: Revising the revisionists." *Review of Economics and Statistics* 90 (2008): 300-323.

Autor, David H., Lawrence F. Katz, and Alan Krueger. "Computing inequality: Have computers changed the labor market?" *Quarterly Journal of Economics* 113 (1998): 1169-1213.

Autor, David H., Frank Levy, and Richard J. Murnane. "The skill content of recent technological change: An empirical exploration." *Quarterly Journal of Economics* 118 (2003): 1279-1334.

Baker, Scott R., Nicholas Bloom, and Steven J. Davis. "Measuring economic policy uncertainty." *Quarterly Journal of Economics* 131 (2016): 1593-1636.

Bekkerman, Ron, and James Allan. "Using bigrams in text categorization." Technical report IR-408, Center of Intelligent Information Retrieval, University of Massachusetts at Amherst (2004).

Berman, Eli, John Bound, and Zvi Griliches. "Changes in the demand for skilled labor within U.S. manufacturing: Evidence from the Annual Survey of Manufacturers." *Quarterly Journal of Economics* 109 (1994): 367-397.

Bloom, Nicholas, Rafaella Sadun, and John Van Reenen. "Management as a technology?" Working paper no. 22327, National Bureau of Economic Research (2016).

Buera, Francisco J., and Ezra Oberfield. "The global diffusion of ideas." *Econometrica* 88 (2020): 83-114.

Bushee, Brian J., Dawn A. Matsumoto, and Gregory S. Miller. "Open versus closed conference calls: The determinants and effects of broadening access to disclosure." *Journal of Accounting and Economics* 34 (2003): 149-180.

Bybee, Leland, Bryan T. Kelly, Asaf Manela, and Dacheng Xiu. "The structure of economic news." Working paper no. 26648, National Bureau of Economic Research (2020).

Caprettini, Bruno, and Hans-Joachim Voth. "Rage against the machines: Labor-saving technology and unrest in industrializing England." *American Economic Review: Insights* 2 (2020): 305-320.

Caselli, Francesco. "Technological revolutions." *American Economic Review* 89 (1999): 78-102.

Caselli, Francesco, and Wilbur John Coleman. "Cross-country technology diffusion: The case of computers." *American Economic Review*, 91 (2001): 328-35

Christensen, Clayton M. *The Innovator's Dilemma: When New Technologies Cause Great Firms to Fail*. Boston, Harvard Business Review Press (1997).

Comin, Diego, and Bart Hobijn. "Cross-country technology adoption: Making the theories face the facts." *Journal of Monetary Economics* 51 (2004): 39-83.

Comin, Diego, and Bart Hobijn. "An exploration of technology diffusion." *American Economic Review* 100 (2010): 2031–59.

Davies, Mark. "The 385+ million word Corpus of Contemporary American English (1990–2008+): Design, architecture, and linguistic insights." *International Journal of Corpus Linguistics* 14 (2009): 159-190.

Deming, David J., and Kadeem Noray. "Earnings dynamics, changing job skills, and STEM careers." *Quarterly Journal of Economics*. 135 (2020): 1965-2005.

Flynn, Joel P., and Karthik A. Sastry. "The macroeconomics of narratives," Working paper 4140751, SSRN (2022).

Forman, Chris, Avi Goldfarb, and Shane Greenstein, "Agglomeration of invention in the Bay Area: Not just ICT." *American Economic Review Papers and Proceedings* 106 (2016): 146-151.

Furman, Jeffrey L., and Megan J. MacGarvie. "Academic science and the birth of industrial research laboratories in the U.S. pharmaceutical industry." *Journal of Economic Behavior and Organization* 63 (2007): 756-776.

Galor, Oded, and Omer Moav. "Ability-biased technological transition, wage inequality, and economic growth." *Quarterly Journal of Economics* 115 (2000): 469-497.

Gentzkow, Matthew, Bryan Kelly, and Matt Taddy. "Text as data." *Journal of Economic Literature* 57 (2019): 535-574.

Glaeser, Edward L., Sari P. Kerr, and William R. Kerr. "Entrepreneurship and urban growth: An empirical assessment with historical mines." *Review of Economics and Statistics* 97 (2015): 498-520.

Goldin, Claudia D., and Lawrence F. Katz. "The origins of technology-skill complementarity." *Quarterly Journal of Economics* 113 (1998): 683-732.

Goldin, Claudia D., and Lawrence F. Katz. *The Race Between Education and Technology*. Cambridge, Harvard University Press (2009).

Gompers, Paul, Josh Lerner, and David Scharfstein. "Entrepreneurial spawning: Public corporations and the genesis of new ventures, 1986 to 1999." *Journal of Finance* 60 (2005): 577-614.

Gordon, Robert J. *The Rise and Fall of American Growth: The U.S. Standard of Living since the Civil War*. Princeton, Princeton University Press (2016).

Greenstone, Michael, Richard Hornbeck and Enrico Moretti, "Identifying agglomeration spillovers: Evidence from winners and losers of large plant openings." *Journal of Political Economy* 118 (2010): 536-598.

Griliches, Zvi, "Hybrid corn: An exploration in the economics of technological change." *Econometrica* 25 (1957): 501–22.

Hall, Bronwyn H., Jacques Mairesse, and Laure Turner. "Identifying age, cohort, and period effects in scientific research productivity: Discussion and illustration using simulated and actual data on French physicists." *Economics of Innovation and New Technologies* 16 (2007): 159-177.

Handley, Kyle, and J.F. Li. "Measuring the effects of firm uncertainty on economic activity: New evidence from one million documents." Working paper no. 27896, National Bureau of Economic Research (2020).

Hassan, Tarek A., Stephan Hollander, Laurence van Lent, and Ahmed Tahoun. "Firm-level political risk: Measurement and effects." *Quarterly Journal of Economics* 134 (2019): 2135–2202.

Hassan, Tarek A., Stephan Hollander, Laurence van Lent, and Ahmed Tahoun. "Firm-level exposure to epidemic diseases: Covid-19, SARS, and H1N1." Working paper no. 26971, National Bureau of Economic Research (2021).

Hershbein, Brad, and Lisa B. Kahn. "Do recessions accelerate routine-biased technological change? Evidence from vacancy postings." *American Economic Review* 108 (2018): 1737-72.

Hoberg, Gerard, and Gordon Phillips. "Text-based network industries and endogenous product differentiation." *Journal of Political Economy* 124 (2016): 1423-65.

Jaffe, Adam B. "Real effects of academic research." *American Economic Review* 79 (1989): 957–970.



Jaffe, Adam B., Manuel Trajtenberg, and Rebecca Henderson. "Geographic localization of knowledge spillovers as evidenced by patent citations." *Quarterly Journal of Economics* 108 (1993): 577–598.

Jewkes, John, David Sawers, and Richard Stillerman. *The Sources of Invention*. New York, St. Martin's Press (1969).

Katz, Lawrence F., and Kevin M. Murphy. "Changes in relative wages, 1963-1987: Supply and demand factors." *Quarterly Journal of Economics* 107 (1992): 35-78.

Kelly, Bryan, Dimitris Papanikolaou, Amit Seru, and Matt Taddy. "Measuring technological innovation over the long run." *American Economic Review: Insights* 3 (2021): 303-320.

Kogan, Leonid, Dimitris Papanikolaou, Lawrence D. Schmidt, and Bryan Seegmiller. "Technology, vintage-specific human capital, and labor displacement: Evidence from linking patents with occupations." Working paper no. 29552, National Bureau of Economic Research (2022).

Krueger, Alan B. "How computers have changed the wage structure: Evidence from microdata, 1984-1989." *Quarterly Journal of Economics* 108 (1993): 33-60.

Krusell, Per, Lee E. Ohanian, José-Víctor Ríos-Rull, and Giovanni L. Violante. "Capital-skill complementarity and inequality: A macroeconomic analysis." *Econometrica* 68 (2000): 1029-53.

Lanjouw, Jean O., Ariel Pakes, and Jonathan Putnam. "How to count patents and value intellectual property: The uses of patent renewal and application data." *Journal of Industrial Economics* 46 (1998): 405-432.

Lerner, Josh, and Amit Seru. "The use and misuse of patent data: Issues for finance and beyond." *Review of Financial Studies* 35 (2022): 2667–2704.

Lind, Nelson, and Natalia Ramondo. (2022). "Global innovation and knowledge diffusion." Working paper no. 29629, National Bureau of Economic Research (2022).

Matsumoto, Dawn, Maarten Pronk, and Erik Roelofsen. "What makes conference calls useful? The information content of managers' presentations and analysts' discussion sessions." *Accounting Review* 86 (2011): 1383-1414.

Merton, Robert K. "The Matthew effect in science: The reward and communication systems of science are considered." *Science* 159 (3810) (1969): 56-63.

Michaels, Guy, Ashwini Natraj, and John Van Reenen. "Has ICT polarized skill demand? Evidence from eleven countries over 25 years." *Review of Economics and Statistics* 96 (2014): 60-77.

Mokyr, Joel. *The Lever of Riches: Technological Creativity and Economic Progress*. New York, Oxford University Press (1992).

Moretti, Enrico. "The effect of high-tech clusters on the productivity of top inventors." *American Economic Review* 111 (2021): 3328-75.

Moscona, Jacob. "Environmental catastrophe and the direction of invention: Evidence from the American Dust Bowl." Unpublished working paper, Massachusetts Institute of Technology (2020).

Moser, Petra, Alessandra Voena, and Fabian Waldinger. "German Jewish émigrés and US invention." *American Economic Review* 104 (2014): 3222-55.

Nelson, Richard R., and Edmund S. Phelps. "Investment in humans, technological diffusion, and economic growth." *American Economic Review* 56 (1966): 69-75.

Organisation for Economic Cooperation and Development. *The Measurement of Scientific and Technological Activities: Proposed Guidelines for Collecting and Interpreting Technological Innovation Data* (third edition). Paris, OECD, Chapter 3 (2005).

Piketty, Thomas, and Emmanuel Saez. "Income inequality in the United States, 1913–1998." *Quarterly Journal of Economics* 118 (2003): 1-41.

Popp, David. "Induced innovation and energy prices." *American Economic Review* 92 (2002): 160-180.

Rogers, Everett M., *Diffusion of Innovations*. New York, Free Press (1962).

Sautner, Zacharias, Laurence Van Lent, Grigory Vilkov, and Ruishen Zhang. "Firm-level climate change exposure." *Journal of Finance* 78 (2023): 1449-98.

Schmookler, Jacob. *Invention and Economic Growth*. Cambridge, Harvard University Press (1966).

Schumpeter, Joseph A. *Capitalism, Socialism, and Democracy*. New York, Harper (1942).

Song, Jae, David J. Price, Faith Guvenen, Nicholas Bloom, and Till Von Wachter. "Firming up inequality." *Quarterly Journal of Economics* 134 (2019): 1-50.

Squicciarini, Mara P., and Nico Voigtländer. "Human capital and industrialization: Evidence from the age of enlightenment." *Quarterly Journal of Economics* 130 (2015): 1825-83.

Syverson, Chad. "What determines productivity?" *Journal of Economic Literature* 49 (2011): 326-365.

Tambe, Prasanna. "Big data investment, skills, and firm value." *Management Science* 60 (2014): 1452-69.

Tambe, Prasanna, and Lorin M. Hitt. "Now IT's personal: Offshoring and the shifting skill composition of the US information technology workforce." *Management Science* 58 (2012): 678-695.

Tambe, Prasanna, Lorin Hitt, Daniel Rock, and Erik Brynjolfsson. "Digital capital and superstar firms." Working paper no. 28285, National Bureau of Economic Research (2020).

Tan, Chade-Meng, Yuan-Fang Wang, and Chan-Do Lee. "The use of bigrams to enhance text categorization." *Information Processing & Management* 38 (2002): 529–546.

Tyson Laura D., and Michael Spence. "Exploring the effects of technology on income and wealth inequality." In Heather Boushey, J. Bradford DeLong, and Marshall Steinbaum (editors), *After Piketty: The Agenda for Economics and Inequality*. Cambridge, Harvard University Press, pp. 170–208 (2017).

United States Patent and Trademark Office. *Performance and Accountability Report*. Washington, USPTO (2020).

Vance, J.D. "One on one interview with Ohio US Senate candidate JD Vance," <https://www.youtube.com/watch?v=4FIapZ88BJQ> (October 19, 2022).

Webb, Michael. "The impact of artificial intelligence on the labor market." Unpublished working paper, Stanford University (2020).

Zucker, Lynne, Michael Darby, and Marilyn B. Brewer. "Intellectual human capital and the birth of U.S. biotechnology enterprises." *American Economic Review* 88 (1998): 290-306.

Table 1 – Top technologies by year

Emergence year	Wikipedia title (technology)	Technology bigrams
1979	Hard disk drive	hard disk; disk drive
1980	Barcode reader	barcode reader; code reader; code scanner; barcode scanner
1981	Laser diode	emitting laser; diode laser; semiconductor laser; laser diode
1982	Personal computer	personal computer
1983	Flat-panel display	panel display; flat panel
1984	User interface	user interface
1985	Mobile phone	mobile telephone; cellular telephone; phones mobile; cellular phone; mobile phone; cell phone
1986	Facial recognition system	frt system; recognition software; recognition system; recognition technology; facial recognition
1987	Digital video	digital video
1988	Model organism	animal model
1989	Mobile device	held computer; computer device; handheld computer; mobile device
1990	Debit card	cards debit; card debit; debit card
1991	Flash memory	flash device; nand flash; flash memory
1992	Machine learning	learning algorithm; machine learning
1993	Financial instrument	financial instrument
1994	Active users	active user
1995	Hybrid electric vehicle	hybrid electric
1996	Digital content	digital content
1997	Multicore processor	multi core; core processor
1998	Information privacy	data protection
1999	Unmanned aerial vehicle	aerial vehicle; unmanned aerial
2000	Transaction account	transaction account
2001	Smartphone	smart phone
2002	Online game	online game
2003	Social networking service	networking site; social networking
2004	Electronic discovery	electronic format
2005	LED circuit	led driver
2006	Augmented reality	augmented reality
2007	Self-driving car	autonomous vehicle

**Notes:** This table reports the top technology by number of mentions in earnings calls (in column 2) for every year of emergence between 1976 and 2007 (in column 1). Column 3 lists the associated technology bigram(s). For the year of emergence 1999, the most frequent technology in earnings calls was “adverse event.” We replace “adverse event” (as it gets dropped in our human audit) with the next most frequent technology, “unmanned aerial vehicle.” See Section 2.c of the main text for details.

Table 2 – Examples of technologies and pioneer locations

Machine Learning (1992)			Digital Imaging (1992)		
CBSA	State	Pct. Patents	CBSA	State	Pct. Patents
New York-Newark-Jersey City	NY-NJ-PA	24%	Rochester	NY	18%
Seattle-Tacoma-Bellevue	WA	13%	San Jose-Sunnyvale-Santa Clara	CA	12%
San Jose-Sunnyvale-Santa Clara	CA	12%	San Francisco-Oakland-Hayward	CA	7%
San Francisco-Oakland-Hayward	CA	9%	Fort Collins	CO	6%
			Greeley	CO	5%
			Worcester	MA-CT	4%
Hybrid Electric (1995)			Smart Phone (2001)		
CBSA	State	Pct. Patents	CBSA	State	Pct. Patents
Detroit-Warren-Dearborn	MI	33%	San Francisco-Oakland-Hayward	CA	18%
Ann Arbor	MI	10%	San Jose-Sunnyvale-Santa Clara	CA	18%
Indianapolis-Carmel-Anderson	IN	8%	Seattle-Tacoma-Bellevue	WA	6%
			New York-Newark-Jersey City	NY-NJ-PA	5%
			Los Angeles-Long Beach-Anaheim	CA	4%

**Notes:** The table shows pioneer CBSAs (in column 1), along with their state (in column 2) and the percentage of early cite-weighted patents accounted for by these CBSAs (in column 3) for a sample of four example technology bigrams – “machine learning,” “digital imaging,” “hybrid electric,” and “smart phone.” Early patents are defined as patents filed within ten years of the emergence year of technology. Each technology bigram’s emergence year is given in parentheses. See Section 2.c of the main text for details.

Table 3 – Geographic concentration of patents, skill, and employment

	Total Number	Share Top 5 CBSAs	Top 5 CBSAs
	(1)	(2)	(3)
Panel A: Geographic concentration of U.S. patents			
Disruptive	1,044,351	42.1%	<b>San Jose-Sunnyvale-Santa Clara CA</b> <b>San Francisco-Oakland-Hayward, CA</b> <b>New York-Newark-Jersey City, NY-NJ-PA</b> <b>Seattle-Tacoma-Bellevue, WA</b> <b>Boston-Cambridge-Newton, MA-NH</b>
Intermediate	364,054	34.0%	<b>San Jose-Sunnyvale-Santa Clara CA</b> <b>New York-Newark-Jersey City, NY-NJ-PA</b> <b>San Francisco-Oakland-Hayward, CA</b> Los Angeles-Long Beach-Anaheim, CA <b>Boston-Cambridge-Newton, MA-NH</b>
Non-Disruptive	215,395	30.4%	<b>New York-Newark-Jersey City, NY-NJ-PA</b> <b>San Jose-Sunnyvale-Santa Clara CA</b> <b>San Francisco-Oakland-Hayward, CA</b> Philadelphia-Camden-Wilmington, PA-NJ-DE-MD Chicago-Naperville-Elgin, IL-IN-WI
All New Technologies	1,623,800	33.3%	<b>San Jose-Sunnyvale-Santa Clara, CA</b> <b>San Francisco-Oakland-Hayward, CA</b> <b>New York-Newark-Jersey City, NY-NJ-PA</b> Los Angeles-Long Beach-Anaheim, CA <b>Boston-Cambridge-Newton, MA-NH</b>
All Patents	3,146,114	32.4%	<b>San Jose-Sunnyvale-Santa Clara CA</b> <b>New York-Newark-Jersey City, NY-NJ-PA</b> <b>San Francisco-Oakland-Hayward, CA</b> Los Angeles-Long Beach-Anaheim, CA Chicago-Naperville-Elgin, IL-IN-WI
Most Cited	1,044,351	32.7%	<b>San Jose-Sunnyvale-Santa Clara CA</b> <b>San Francisco-Oakland-Hayward, CA</b> <b>New York-Newark-Jersey City, NY-NJ-PA</b> Los Angeles-Long Beach-Anaheim, CA Chicago-Naperville-Elgin, IL-IN-WI
Panel B: Geographic concentration of skill and employment			
College Graduates	51.5 million	22.5%	<b>New York-Newark-Jersey City, NY-NJ-PA</b> Los Angeles-Long Beach-Anaheim, CA Chicago-Naperville-Elgin, IL-IN-WI Washington-Arlington-Alexandria, DC-VA-MD-WV <b>San Francisco-Oakland-Hayward, CA</b>
Employed	156.5 million	18.9%	<b>New York-Newark-Jersey City, NY-NJ-PA</b> Los Angeles-Long Beach-Anaheim, CA Chicago-Naperville-Elgin, IL-IN-WI Dallas-Fort Worth-Arlington, TX Houston-The Woodlands-Sugar Land, TX

**Notes:** This table reports the concentration of patents, skill, and employment across CBSAs in the U.S. The measures of skill and employment are obtained from the 2015 American Communities Survey. A patent is considered disruptive/intermediate/non-disruptive if it mentions at least one bigram associated with disruptive/intermediate/non-disruptive technologies more than once. The row “Most Cited” shows the geographic concentration of the 1,044,351 patents with the most normalized citations for comparison. This number is chosen to equal the number of patents mentioning a disruptive technology. CBSAs in **bold** are those in the top five for disruptive patents.

Table 4 – Region broadening

Panel A: Main specifications			
	<i>Coefficient of Variation</i> $_{\tau,t}$		
Sample	EC $\geq$ 100	All	
	(1)	(2)	(3)
<i>Years since emergence</i> $_{\tau,t}$	-0.068*** (0.026)	-0.065*** (0.024)	-0.153*** (0.012)
Constant	5.577*** (0.645)	6.212*** (0.585)	8.269*** (0.271)
R-squared	0.019	0.013	0.825
N	4,270	8,347	8,347
Bigrams	428	835	835
Bigram FE	NO	NO	YES
Std. Errors (cluster)	Wiki Title	Wiki Title	Wiki Title
<b>Years to zero CV</b>	82.12	95.52	54.07
Panel B: Alternative measures of geographic concentration			
	$\frac{N_{c,\tau,t}^{Top-5}}{N_{c,\tau,t}^{All}}$	<i>Pct.</i> ( $N_{c,\tau,t} \leq 0.1$ )	$\sum_i (N_{c,\tau,t} - 1)^2$
	(1)	(2)	(3)
<i>Years since emergence</i> $_{\tau,t}$	-1.724*** (0.136)	-1.117*** (0.088)	-173.965** (77.247)
Constant	88.781*** (3.181)	95.514*** (2.063)	13,650.448*** (1,808.008)
R-squared	0.846	0.922	0.679
N	8,347	8,347	8,347
Bigrams	835	835	835
Bigram FE	YES	YES	YES
Std. Errors (cluster)	Wiki Title	Wiki Title	Wiki Title
<b>Years to zero CV</b>	51.49	85.48	78.47

**Notes:** This table reports the results from regressions at the technology bigram x year level. The dependent variable is a measure of the geographic concentration of a given technology bigram's job postings in a given year. The independent variable – years since emergence – is the number of years that have elapsed since the technology's year of emergence. Panel A reports results using our baseline measure of geographic concentration – the coefficient of variation of the normalized share of a technology bigram's job postings across CBSAs. Panel B reports results using three alternative measures of geographic concentration – the mean normalized share of a technology's job postings in the top five CBSAs relative to the mean normalized share across all CBSAs, the percentage of CBSAs with a normalized share of a technology's job postings of less than 10% (that is, the representation of CBSAs with almost no activity associated with that bigram), and the sum of squared deviations of the normalized share from 1 (similar to the Herfindahl-Hirschman Index). Column 1 of Panel A is restricted to the sample of technology bigrams that appear in at least 100 earnings calls. The other regressions use all technology bigrams that appear in at least 1000 job postings in our sample. Observations are weighted by the square root of the total number job postings mentioning that technology in that year, capped at 100. The normalized share of job postings is capped at the 99<sup>th</sup> percentile of non-zero observations. Standard errors are clustered by Wikipedia title (technology). All specifications indicate fixed effects used. Years to zero CV are calculated by dividing the constant by the coefficient estimate on the years since emergence.

Table 5 – Pioneer location advantage in technology hiring

Sample:	<i>Normalized Share</i> <sub><i>c,t</i></sub>			
	<i>EC</i> <sub><i>τ</i></sub> ≥ 100		All	
	(1)	(2)	(3)	(4)
<i>Pioneer</i> <sub><i>c,τ</i></sub>	0.311*** (0.076)	1.084*** (0.309)	1.321*** (0.254)	1.282*** (0.243)
<i>Pioneer</i> <sub><i>c,τ</i></sub> * <i>Years since emg</i> <sub><i>τ,t</i></sub>		-0.032** (0.013)	-0.035*** (0.011)	-0.034*** (0.010)
<i>Pioneer Neighbor</i> <sub><i>c,τ</i></sub>				0.158*** (0.057)
<i>Pioneer Neighbor</i> <sub><i>c,τ</i></sub> * <i>Years since emg</i> <sub><i>τ,t</i></sub>				-0.004 (0.003)
R-squared	0.038	0.038	0.030	0.030
N	3,965,122	3,965,122	7,751,024	7,751,024
Bigrams	428	428	835	835
Bigram FE	YES	YES	YES	YES
CBSA FE	YES	YES	YES	YES
Year FE	YES	YES	YES	YES
Std. Errors (cluster)	Wiki Title	Wiki Title	Wiki Title	Wiki Title
<i>Rate of decline per year</i>		-0.029 (0.005)	-0.027 (0.003)	-0.026 (0.003)
<b>Implied years to zero advantage</b>		<b>33.88</b>	<b>37.74</b>	<b>38.26</b>

**Notes:** This table reports results from regressions of the *Normalized Share*<sub>*c,t*</sub> (for each CBSA x technology bigram x year) on a dummy indicating the pioneer status of the CBSA and the interaction of this dummy with the number of years that have elapsed since the bigram's emergence. The dummy variable *Pioneer Neighbor*<sub>*c,τ*</sub> takes value one for non-pioneer CBSAs that are within 100 miles of the technology's pioneer locations. Columns 1 and 2 are restricted to the sample of technology bigrams that appear in at least 100 earnings calls. The other regressions use all technology bigrams that appear in at least 1000 job postings in our sample. Observations are weighted by the square root of the total number job postings mentioning that technology in that year, capped at 100. The normalized share of job postings is capped at the 99<sup>th</sup> percentile of non-zero observations. All specifications indicate fixed effects used. Standard errors are clustered by Wikipedia title (technology). The rate of decline per year is calculated as  $\frac{\beta_D}{\beta_P}$ , where  $\beta_P$  is the coefficient on *Pioneer*<sub>*c,τ*</sub> and  $\beta_D$  is the coefficient of *Pioneer*<sub>*c,τ*</sub> \* *Years since emg*<sub>*τ,t*</sub>.



Table 6 – Mechanisms: Spread of high vs. low-skill jobs;  
Spread of research, development, and production jobs vs. use jobs

Panel A: Region-broadening regressions				
	$\log(\text{Coefficient of Variation})_{\tau,t}$			
	EC $\geq$ 100			
Sample:	(1)	(2)	(3)	(4)
	High-Skill Job Postings	Low-Skill Job Postings	RDP Job Postings	Use Job Postings
$Years\ since\ emergence_{\tau,t}$	-0.030*** (0.003)	-0.042*** (0.004)	-0.018*** (0.003)	-0.035*** (0.003)
R-squared	0.847	0.848	0.804	0.904
N	4,188	4,188	3,380	3,380
Bigram FE	YES	YES	YES	YES
Std. Errors (cluster)	Wiki Title	Wiki Title	Wiki Title	Wiki Title
Panel B: Pioneer advantage regressions				
	$Normalized\ Share_{c,\tau,t}$			
	EC $\geq$ 100			
Sample:	(1)	(2)	(3)	(4)
	High-Skill Job Postings	Low-Skill Job Postings	RDP Job Postings	Use Job Postings
$Pioneer_{c,\tau}$	1.074*** (0.279)	0.939*** (0.282)	2.138*** (0.650)	1.472*** (0.324)
$Pioneer_{c,\tau} * Years\ since\ emg_{\tau,t}$	-0.026** (0.012)	-0.034*** (0.012)	-0.046 (0.032)	-0.038*** (0.014)
R-squared	0.019	0.014	0.007	0.046
N	4,348,921	4,251,374	3,965,122	3,965,122
Bigram FE	YES	YES	YES	YES
CBSA FE	YES	YES	YES	YES
Year FE	YES	YES	YES	YES
Std. Errors (cluster)	Wiki Title	Wiki Title	Wiki Title	Wiki Title
$Rate\ of\ decline\ per\ year$	-0.024 0.006	-0.036 0.003	-0.022 0.009	-0.026 0.004
<b>Implied years to zero advantage</b>	<b>41.30</b>	<b>27.62</b>	<b>46.48</b>	<b>38.74</b>

**Notes:** This table reports region-broadening regressions at the technology x year level (Panel A) and pioneer advantage regressions at the technology bigram x year x CBSA level (Panel B). Columns 1 and 2 show separate regressions for high-skill (column 1) and low-skill (column 2) job postings. Column 3 shows regressions for research, development, and production-related job postings (RDP); column 4 for job postings relating to the use of the technology. For definitions of these concepts, see Section 4.c of the main text. All specifications use technology bigrams that appear in at least 100 earnings calls. Observations are weighted by the square root of the total number job postings mentioning that technology in that year, capped at 100. All specifications indicate fixed effects used. Standard errors are clustered by Wikipedia title (technology). Panel A reports results from regressions of  $\log(\text{Coefficient of Variation})_{\tau,t}$  on the number of years since the technology's year of emergence. In a stacked specification, the difference between coefficient on  $Years\ since\ emergence_{\tau,t}$  in columns 1 and 2 is -0.013 (S.E. = 0.004, p-val. = 0.001). The difference between the coefficient on  $Years\ since\ emergence_{\tau,t}$  in columns 3 and 4 is -0.017 (S.E. = 0.003, p-val = 0.000). Both differences are thus statistically distinguishable from zero. Panel B reports results from regressions of the  $Normalized\ share_{c,\tau,t}$  (for each CBSA, bigram, and year) on a dummy indicating pioneer status of the CBSA and on the interaction of this dummy with the number of years that have elapsed since bigram's emergence. The normalized share of job postings is capped at the 99<sup>th</sup> percentile of non-zero observations. In a stacked specification, the difference between estimates of Rate of decline per year in columns 1 and 2 is 0.011 (S.E. = 0.005, p-val = 0.026), while the difference between  $\beta_D$  (coefficient of  $Pioneer_{c,\tau} \times Years\ since\ emergence_{\tau,t}$ ) is 0.022 (S.E. = 0.011, p-val = 0.056). Similarly, the difference between Rate of decline per year in columns 3 and 4 is 0.004 (S.E. = 0.007, p-val = 0.518), while the difference between  $\beta_D$  (coefficient of  $Pioneer_{c,\tau} \times Years\ since\ emergence_{\tau,t}$ ) is 0.008 (S.E. = 0.026, p-val = 0.752). The rate of decline per year is calculated as  $\frac{\beta_D}{\beta_P}$ , where  $\beta_P$  is the coefficient on  $Pioneer_{c,\tau}$  and  $\beta_D$  is the coefficient of  $Pioneer_{c,\tau} * Years\ since\ emg_{\tau,t}$ .

Table 7 – Skill broadening

Panel A: Main specifications				
Sample	<i>Share College Educated<sub>τ,t</sub> * 100</i>			
	EC >= 100	All		
	(1)	(2)	(3)	(4)
Constant	57.078*** (2.135)	59.095*** (1.794)	57.475*** (2.294)	63.898*** (0.840)
<i>Years since emergence<sub>τ,t</sub></i>	-0.228** (0.092)	-0.288*** (0.079)	-0.218*** (0.100)	-0.493*** (0.036)
R-squared	0.017	0.019	0.024	0.910
N	4,270	8,347	8,347	8,347
Bigrams	428	835	835	835
Year FE	NO	NO	YES	NO
Bigram FE	NO	NO	NO	YES
Standard Errors (cluster)	Wiki Title	Wiki Title	Wiki Title	Wiki Title
<b>Implied years to average skill</b>	<b>117.23</b>	<b>100.03</b>	<b>124.37</b>	<b>68.08</b>
Panel B: Alternative measures of skill				
	(1)	(2)	(3)	
	<i>Years of Schooling<sub>τ,t</sub></i>	<i>Share Post Graduates<sub>τ,t</sub> * 100</i>	<i>Average Wage<sub>τ,t</sub></i>	
Constant	15.504*** (0.047)	22.617*** (0.456)	75,521.317*** (840.562)	
<i>Years since emergence<sub>τ,t</sub></i>	-0.024*** (0.002)	-0.149*** (0.020)	-505.134*** (35.986)	
R-squared	0.915	0.905	0.889	
N	8,347	8,347	8,347	
Bigrams	835	835	835	
Bigram FE	YES	YES	YES	
Std. Errors (cluster)	Wiki Title	Wiki Title	Wiki Title	
<b>Implied years to average skill</b>	<b>77.59</b>	<b>78.03</b>	<b>69.72</b>	

**Notes:** This table reports results from regressions at the technology bigram x year level. The dependent variable is a measure of the average skill requirement of a technology bigram’s job postings in a given year. The independent variable is the number of years that have elapsed since the technology’s emergence. The dependent variable in Panel A is the average share of job postings mentioning technology bigram  $\tau$  in year  $t$  that require a college degree. Panel B shows results corresponding to column 4 of Panel A for alternative measures of skill associated with technology bigram job postings: average years of schooling (column 1), share of post-graduates (in column 2), and average wage (in column 3). Column 1 of Panel A is restricted to the sample of technology bigrams that appear in at least 100 earnings calls. The other regressions use all technology bigrams that appear in at least 1000 job postings in our sample. Observations are weighted by the square root of the total number job postings mentioning that technology in that year, capped at 100. All specifications indicate fixed effects used. Standard errors are clustered by Wikipedia title (technology). The row “Implied years to average skill” is determined by  $-(Constant - Average Population Skill)/\beta_{SB}(Years since emergence_{\tau,t})$ , where *Average Population Skill* represents the weighted average skill of the US population according to the 2015 ACS Survey.

Table 8 – Skill broadening mechanisms: Research, development, and production and training jobs

Sample	<i>Share College Educated<sub>τ,t</sub></i>				
	EC >= 100				
	(1)	(2)	(3)	(4)	(5)
<i>Years since emergence<sub>τ,t</sub></i>	-0.228** (0.092)	-0.193** (0.076)	-0.186** (0.076)	-0.177** (0.075)	-0.232*** (0.081)
<i>Share of R&amp;D Postings<sub>τ,t</sub></i>		1.038*** (0.123)			
<i>Share of Produce Postings<sub>τ,t</sub></i>			0.388*** (0.078)		
<i>Share of RDP Postings<sub>τ,t</sub></i>				0.384*** (0.071)	
<i>Share Training Required<sub>τ,t</sub></i>					0.270*** (0.048)
<i>Constant</i>	57.078*** (2.135)	52.565*** (1.841)	51.803*** (2.009)	51.246*** (1.998)	52.795*** (2.070)
R-squared	0.017	0.206	0.202	0.215	0.136
N	4,270	4,270	4,270	4,270	4,270
Standard Errors (Cluster)	Wiki Title	Wiki Title	Wiki Title	Wiki Title	Wiki Title

**Notes:** This table reports results from regressions at the technology bigram x year level. Column 1 replicates the specification in Table 7, Panel A, column 1. Columns 2-5 add additional controls: the share of the technology's job postings relating to research and development (column 2), the share of the technology's job postings relating to the technology's production (column 3), the share of the technology's job postings relating to research, development, and production (column 4), and the share of the technology's job postings requiring training in the technology (column 5). The regressions are restricted to the sample of technology bigrams that appear in at least 100 earnings calls. Observations are weighted by the square root of the total number job postings mentioning that technology in that year, capped at 100. All specifications indicate fixed effects used. Standard errors are clustered by Wikipedia title (technology).

Table 9 – Broadening and pioneer advantage across different dimensions

Panel A: Broadening				
	<i>Coefficient of Variation<sub>τ,t</sub></i>			
	(1)	(2)	(3)	(4)
	Industries	Occupations	Firms	CBSAs
<i>Years since emergence<sub>τ,t</sub></i>	-0.018 (0.017)	-0.056*** (0.015)	-0.354*** (0.038)	-0.153*** (0.012)
<i>Constant</i>	4.928*** (0.400)	8.136*** (0.351)	22.042*** (0.890)	8.269*** (0.271)
R-squared	0.817	0.763	0.919	0.825
N	4,970	8,347	4,580	8,347
Bigrams	497	835	458	835
Bigram FE	YES	YES	YES	YES
Std. Errors (cluster)	Wiki Title	Wiki Title	Wiki Title	Wiki Title
Mean	4.52	6.83	13.78	4.69
<b><i>Rate of decline per year</i></b>	<b>-0.004</b> <b>(0.003)</b>	<b>-0.007</b> <b>(0.002)</b>	<b>-0.016</b> <b>(0.001)</b>	<b>-0.018</b> <b>(0.001)</b>
Years to zero CV	405.50	279.11	62.21	54.07
Panel B: Pioneer advantage				
	<i>Normalized Share<sub>n,τ,t</sub></i>			
	(1)	(2)	(3)	
	Industries	Firms	CBSAs	
<i>Pioneer<sub>n,τ</sub></i>	6.504*** (1.751)	20.935*** (4.883)	1.321*** (0.254)	
<i>Pioneer<sub>n,τ</sub> * Years since emg<sub>τ,t</sub></i>	-0.082 (0.074)	-0.489*** (0.185)	-0.035*** (0.011)	
R-squared	0.043	0.009	0.030	
N	1,515,850	49,854,895	7,751,024	
Bigrams	497	458	835	
Bigram FE	YES	YES	YES	
CBSA FE	YES	YES	YES	
Year FE	YES	YES	YES	
Std. Errors (cluster)	Wiki Title	Wiki Title	Wiki Title	
<b><i>Rate of decline per year</i></b>	<b>-0.013</b> <b>(0.008)</b>	<b>-0.023</b> <b>(0.004)</b>	<b>-0.027</b> <b>(0.003)</b>	
Implied years to zero advantage	79.32	42.81	37.74	

**Notes:** This table reports results from broadening regressions (in Panel A) and pioneer advantage regressions (in Panel B) along four dimensions: 1) industries, 2) occupations, 3) firms, and 4) locations (CBSAs). In Panel A, we regress coefficient of variation calculated over *Normalized Share<sub>n,τ,t</sub>* for each bigram and year where *n* is an industry (in column 1), occupation (in column 2), firm (in column 3), and location (in column 4). Panel B reports results from regressions of the *Normalized share<sub>n,τ,t</sub>* on the pioneer status of *n* and the interaction of the pioneer status with the year since the technology bigram's emergence. As in Panel A, *n* is an industry (in column 1), firm (in column 2), and location (in column 3). The regressions use all technology bigrams that appear in at least 1000 job postings in our sample. All specifications are weighted by the square root of the total number job postings mentioning that technology in that year, capped at 100. The normalized share is capped at 99th percentile of non-zero observations. All specifications indicate fixed effects used. Standard errors are clustered by Wikipedia title (technology). Note that the number of bigrams changes across specifications, depending on data availability on firms and industries in job postings. To test whether estimated coefficients are different across dimensions, we estimate stacked regressions using the same specifications as in Panels A and B, where we interact fixed effects with indicators for each dimension. In Panel A, the absolute rate of decline across CBSAs is 0.017 (0.005)\*\*\*, 0.012 (0.003)\*\*\*, and 0.004 (0.002)\* higher across than industries, occupations, and firms, respectively. In Panel B, the absolute rate of decline in pioneer advantage across CBSAs is 0.014 (0.007)\* higher than across industries and 0.003 (0.005) higher than across firms. Similarly, the coefficient of *Pioneer<sub>n,τ</sub>* is 19.614 (6.604)\*\*\* and 5.183 (1.659)\*\*\* higher for firms and industries than for CBSAs. For the coefficient on *Pioneer<sub>i,τ</sub> \* Years since emg<sub>τ,t</sub>*, the estimated differences are 0.454 (0.247)\* and 0.047 (0.702) between CBSAs and, respectively, firms and industries.

Table 10 – Robustness checks: Alternative samples and specifications

	Share of Top-5 CBSAs (1)	Coefficient of Variation (2)	log(Coefficient of Variation) (3)	Normalized Share (4)	Share College Educated (5)
	Concentration of Innovation [Table 3, Panel A, row 1 col. 2]	Region Broadening [Table 4, Panel A, col. 3]	Region Broadening by Skill [Table 6, Panel A, col. 1, 2]	Rate of decline in Pioneer Persistence [Table 5, col. 2]	Skill Broadening [Table 7, Panel A, col. 4]
Estimate/Coefficient:	Share of Top-5 CBSAs	$\beta_{RB}$	$\beta_{RB}^{High\ skill} - \beta_{RB}^{Low\ skill}$	$\beta_D/\beta_P$	$\beta_{SB}$
<b>Panel A: Influential patents</b>					
Baseline: At least 1,000 cite-wt. patents	42.1%	-0.153*** (0.012)	-0.013*** (0.004)	-0.029*** (0.005)	-0.493*** (0.036)
At least 1,250 cite-wt. patents	42.2%	-0.150*** (0.012)	-0.012*** (0.004)	-0.029*** (0.005)	-0.488*** (0.037)
At least 1,500 cite-wt. patents	42.4%	-0.146*** (0.012)	-0.012*** (0.004)	-0.029*** (0.005)	-0.500*** (0.038)
At least 1,750 cite-wt. patents	42.5%	-0.146*** (0.013)	-0.011*** (0.004)	-0.029*** (0.005)	-0.501*** (0.040)
At least 2,000 cite-wt. patents	42.5%	-0.147*** (0.013)	-0.011** (0.004)	-0.029*** (0.006)	-0.494*** (0.040)
<b>Panel B: Phrase Length</b>					
Baseline: Bigrams	42.1%	-0.153*** (0.012)	-0.013*** (0.004)	-0.029*** (0.005)	-0.493*** (0.036)
Bigrams and trigrams	42.0%	-0.127*** (0.015)	-0.013*** (0.004)	-0.029*** (0.005)	-0.492*** (0.032)
Bigrams, trigrams, and unigrams (disruptive only)	37.0%	-0.106*** (0.008)	-0.011*** (0.002)	-0.029*** (0.003)	-0.498*** (0.034)
<b>Panel C: Human Audit</b>					
Baseline: Bigrams	42.1%	-0.153*** (0.012)	-0.013*** (0.004)	-0.029*** (0.005)	-0.493*** (0.036)
Human-audited bigrams (disruptive only)	44.3%	-0.142*** (0.019)	-0.015*** (0.005)	-0.025*** (0.009)	-0.574*** (0.061)
<b>Panel D: Alternative emergence years</b>					
Baseline: At least 100 cite-wt. patents	42.1%	-0.153*** (0.012)	-0.013*** (0.004)	-0.029*** (0.005)	-0.493*** (0.036)
At least 1 cite-wt. patent	41.9%	-0.157*** (0.011)	-0.014*** (0.004)	-0.030*** (0.003)	-0.465*** (0.034)
At least 200 cite-wt. patents	42.2%	-0.154*** (0.013)	-0.015*** (0.004)	-0.031*** (0.005)	-0.511*** (0.042)
50% of total cite-wt. patents	39.2%	-0.144*** (0.009)	-0.012*** (0.004)	-0.034*** (0.011)	-0.466*** (0.027)
<b>Panel E: Alternative weighting schemes</b>					
Baseline: min(100, sqrt(# postings))	NA	-0.153*** (0.012)	-0.013*** (0.004)	-0.029*** (0.005)	-0.493*** (0.036)
Unweighted regression	NA	-0.194***	-0.012***	-0.026***	-0.532***

All bigrams with job postings	NA	(0.014) -0.209***	(0.004) -0.012***	(0.004) -0.020***	(0.043) -0.490***
Log-wt: min(100, log(# postings))	NA	(0.014) -0.178***	(0.002) -0.013***	(0.003) -0.027***	(0.046) -0.515***
Bigrams collapsed into technologies	NA	(0.013) -0.156***	(0.004) -0.013***	(0.004) -0.025***	(0.039) -0.490***
		(0.010)	(0.003)	(0.005)	(0.033)

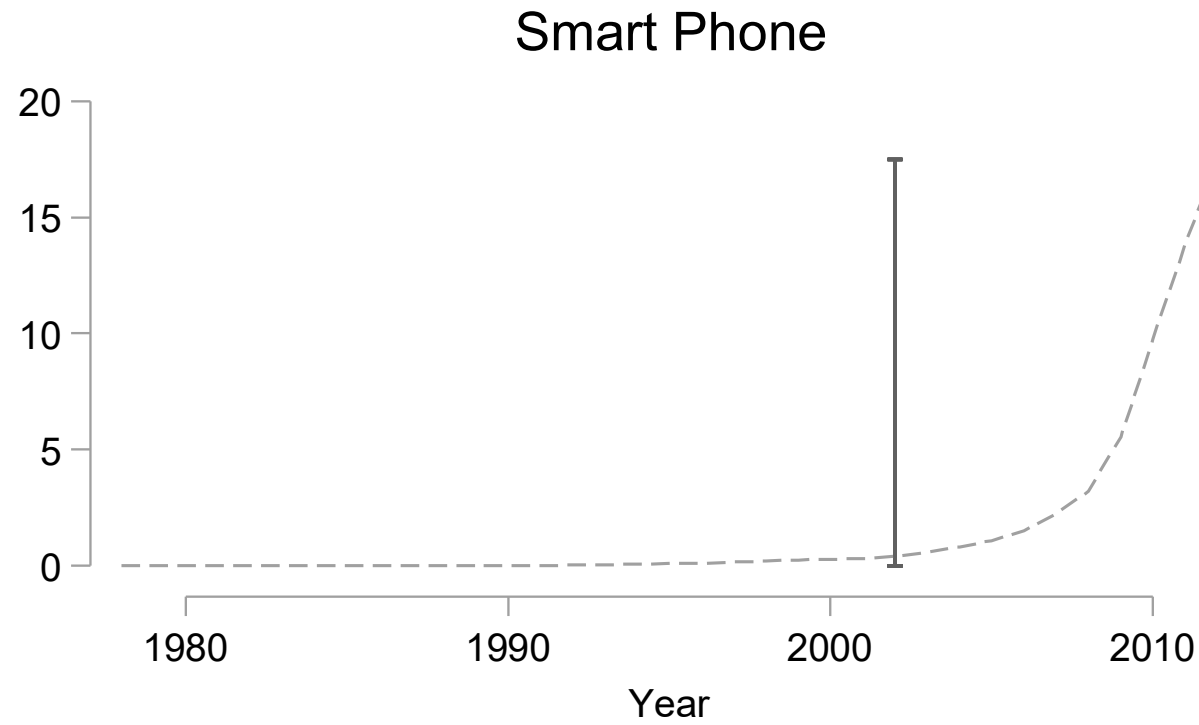
**Notes:** This table reports robustness checks to our primary (“Baseline”) results. Panel A reports robustness to changing the threshold for defining bigrams associated with “influential innovations.” In the baseline, we retain only those that appear in patents accumulating a total of at least 1,000 weighted citations. The panel shows four variations, with cutoffs ranging from 1,250 to 2,000 citations. Panel B presents results from extending our sample to include technology trigrams and unigrams. In the baseline, we include only technology bigrams. The row “Bigrams and trigrams” includes trigrams along with technology bigrams in the analysis, while the row “Bigrams, trigrams, and unigrams (disruptive only)” further adds unigrams. While adding unigrams, we restrict the sample to only disruptive ones (those appearing in 100 or more earning calls). In Panel C, “Human audited bigrams (disruptive only)”, we rely on human reading to determine whether or not a bigram describes a technology, instead of the Wikipedia filter. We report results from including only those technology bigrams that survived the human auditing process and are associated with disruptive technologies. Panel D reports results from variations in defining emergence years. In the baseline, the emergence year is defined as the first year in which (a) 100 citation-weighted patents associated with that technology had been already applied for and (b) where the next five years had 10% annual growth in (smoothed) weighted patenting. The rows “At least 1 cite-wt. patent” and “At least 200 cite-wt. patents” explore changing the threshold from 100 cite-weighted patents in (a) to at least one cite-weighted patent and at least 200 cite-weighted patents, respectively. The row “50% of total cite-wt. patents” changes our definition of emergence years completely: the emergence year of a given bigram is defined as the first year when 50% of maximum peak of citation-weighted patent counts is realized. Panel E presents robustness to changing weighting schemes in regressions. Baseline regressions are weighted by the square root of the total number job postings mentioning that technology in that year, capped at 100. The row “Unweighted regression” replicates baseline regressions with equal weights for each observation. The row “All bigrams with job postings” performs unweighted regressions with all bigrams that are mentioned by at least one job posting. The row “Log-wt” weights observations by the log of the number of postings observed for each technology bigram in a given year, capped at 100. The last row “Bigrams collapsed into technologies” replicates our results when all bigrams associated with a given Wikipedia title (technology) are collapsed into the technology. See the original regressions for full details.

Table 11 – Robustness checks: Alternative specifications of standard errors

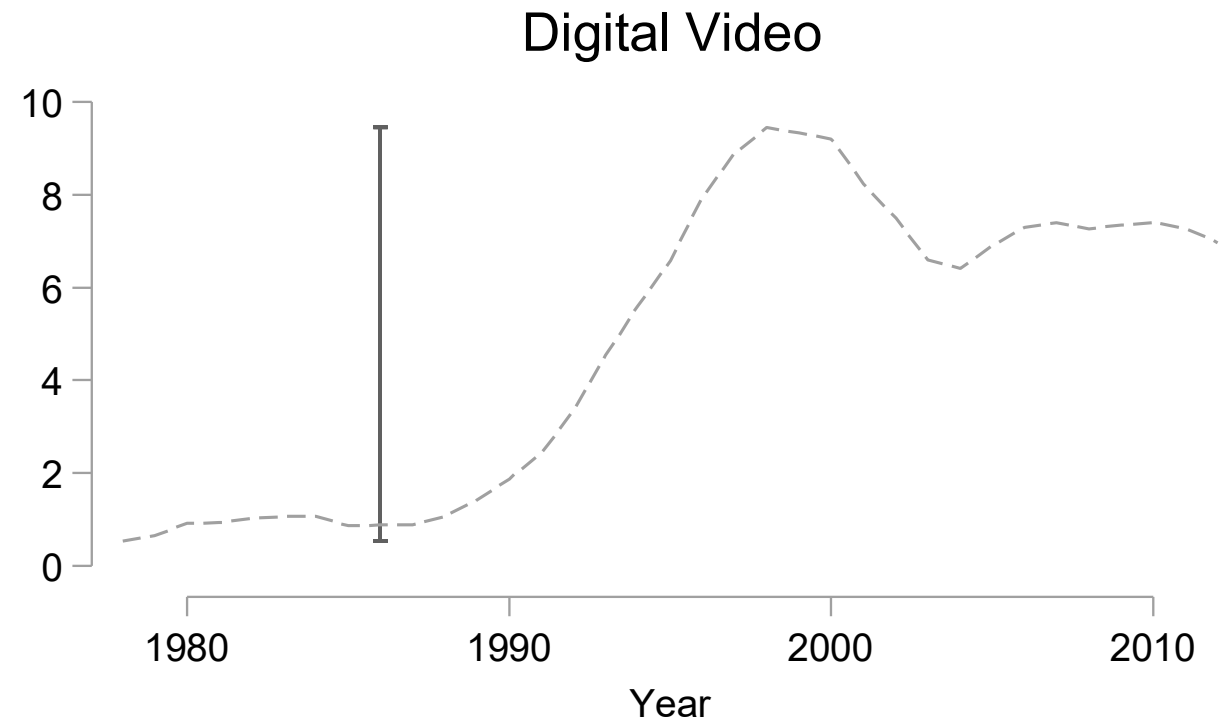
	Coefficient of Variation	Coefficient of Variation	Normalized Share	Share College Educated
	(1)	(2)	(3)	(4)
	Region Broadening	Region Broadening by Skill	Pioneer Persistence	Skill Broadening
	[Table 4, Panel A, col. 3]	[Table 6, Panel A, col 1,2]	[Table 5, col. 2]	[Table 7, Panel A, col. 4]
<i>Years since emergence</i> <sub><math>\tau,t</math></sub> ( <i>High Skill</i> in col. 3)	-0.153	-0.030		-0.493
[baseline] Cluster, Wikipedia Title level	(0.012) ***	(0.003) ***		(0.036) ***
Cluster, Bigram level	(0.009) ***	(0.002) ***		(0.028) ***
Cluster, Year level	(0.016) ***	(0.004) ***		(0.042) ***
Bootstrap (1,000 replications)	(0.009) ***	(0.002) ***		(0.029) ***
<i>Years since emergence</i> <sub><math>\tau,t</math></sub> ( <i>Low Skill</i> )		-0.042		
[baseline] Cluster, Wikipedia Title level		(0.004) ***		
Cluster, Bigram level		(0.003) ***		
Cluster, Year level		(0.003) ***		
Bootstrap (1,000 replications)		(0.003) ***		
<i>Pioneer</i> <sub><math>i,\tau</math></sub>			1.084	
[baseline] Cluster, Wikipedia Title level			(0.309) ***	
Cluster, Bigram level			(0.259) ***	
Cluster, Year level			(0.080) ***	
Cluster, CBSA level			(0.202) ***	
Cluster, State level			(0.201) ***	
Cluster, CBSA-Wikipedia Title levels			(0.309) ***	
Bootstrap (1,000 replications)			(0.254) ***	
<i>Pioneer</i> <sub><math>i,\tau</math></sub> * <i>Years since emergence</i> <sub><math>\tau,t</math></sub>			-0.032	
[baseline] Cluster, Wikipedia Title level			(0.013) **	
Cluster, Bigram level			(0.011) ***	
Cluster, Year level			(0.004) ***	
Cluster, CBSA level			(0.008) ***	
Cluster, CBSA-Wikipedia Title levels			(0.013) **	
Cluster, State level			(0.008) ***	
Bootstrap (1,000 replications)			(0.011) ***	
$\beta(\text{Pioneer}_{i,\tau} * \text{Years since emergence}_{\tau,t})$ / $\beta(\text{Pioneer}_{i,\tau})$			-0.029	
[baseline] Cluster, Wikipedia Title level			(0.005) ***	
Cluster, Bigram level			(0.004) ***	
Cluster, Year level			(0.001) ***	
Cluster, CBSA level			(0.004) ***	
Cluster, CBSA-Wikipedia Title levels			(0.005) ***	
Cluster, State level			(0.003) ***	
Bootstrap (1,000 replications)			(0.004) ***	
Bigram FE	YES	YES	YES	YES
Skill FE	NA	YES	NA	NA
CBSA FE	NA	NA	YES	NA
Year FE	NA	NA	YES	NA

**Notes:** This table reports results from varying specifications for standard errors corresponding to coefficient estimates for our main results – region broadening, region broadening by skill, pioneer persistence and skill broadening. The statistical significance of coefficients is indicated by the asterisks next to each parenthesis. For the results in Columns 1, 2, and 4, we report standard errors clustered at the Wikipedia title level (baseline), bigram level, and year level. In Column 3, we report standard errors clustered at the Wikipedia title level (baseline), bigram level, year level, CBSA level, state level, and CBSA x Wikipedia title level (double-cluster). To cluster CBSAs into the state level, we assign CBSAs that are shared by more than one state to the state with lowest FIPS number. For each result, in the last row, we report bootstrapped standard errors for each specification. Bootstrapped standard errors are computed based on 1,000 replications with replacement from the original sample. Re-sampling was done at the bigram-level (sampling bigram-blocks with ten years of observations). See the original regressions for full details.

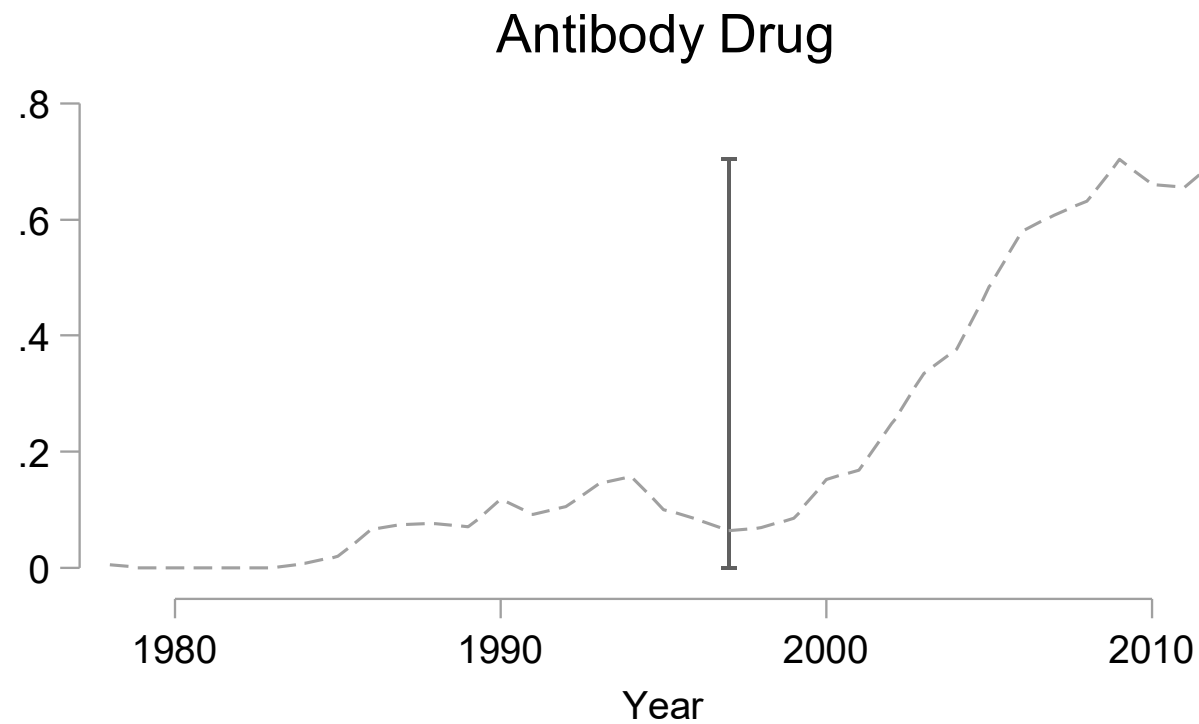
**Figure 1– Examples of emergence year definition**



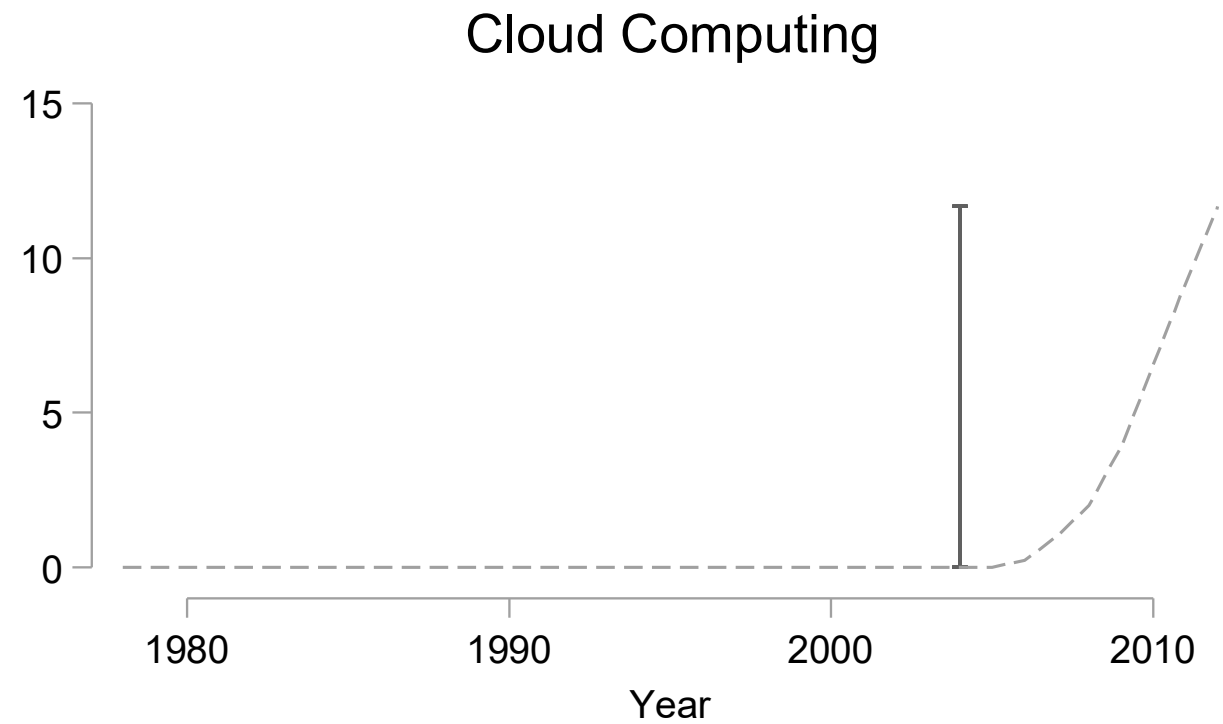
—|— Emergence year    - - - Patent series



—|— Emergence year    - - - Patent series



—|— Emergence year    - - - Patent series

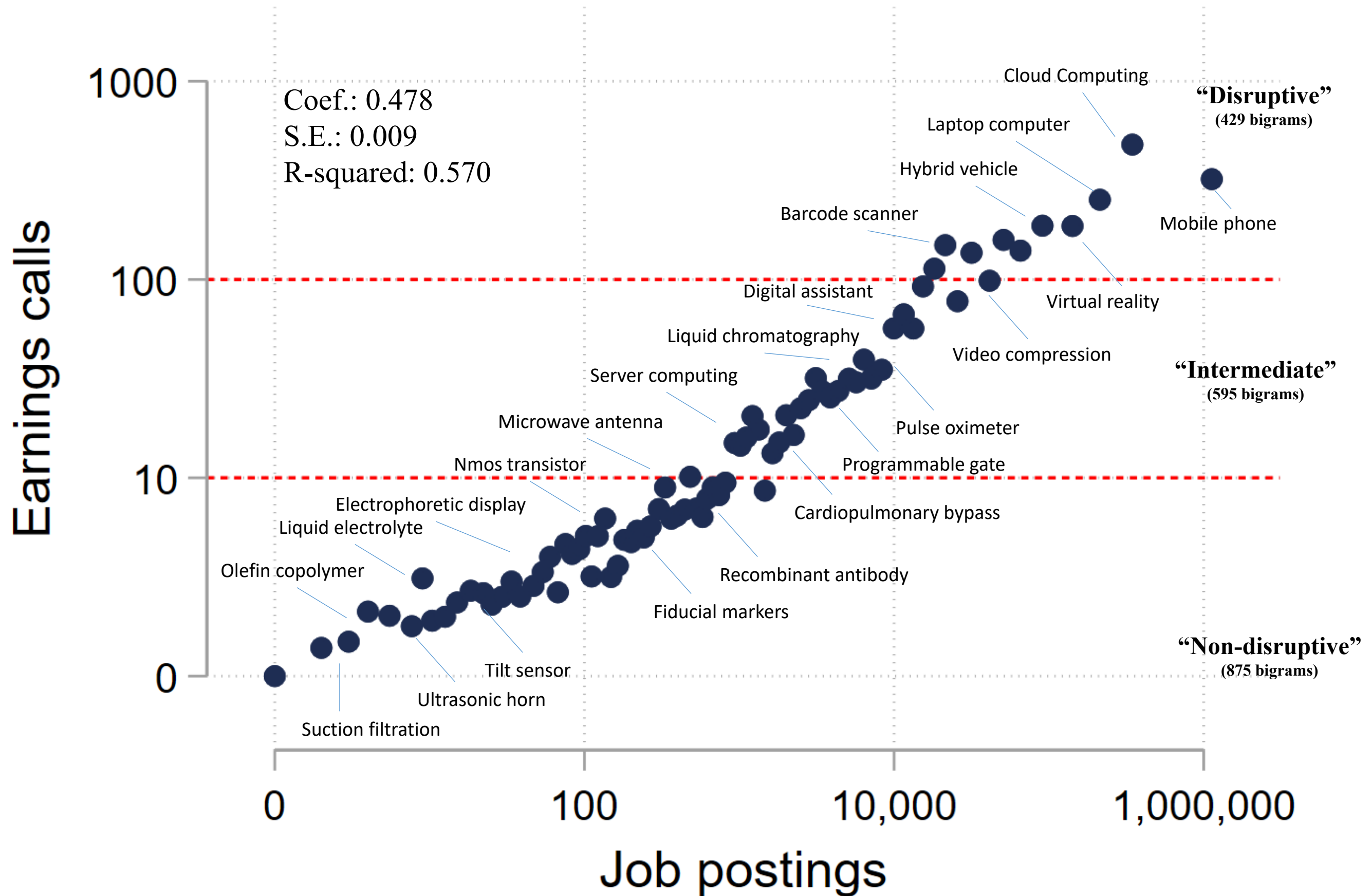


—|— Emergence year    - - - Patent series

**Notes:** The figure shows four examples of attribution of emergence years. In each example, the time series plots the smoothed number of cite-weighted patents associated with the technology by year of application of the patent. For each bigram, we mark the emergence year as the first year in which (a) the technology reaches 100 cite-weighted patent applications and (b) where the next five years had at least 10% annual growth in the (smoothed) series for each bigram. For more details, refer to Section 2.c.



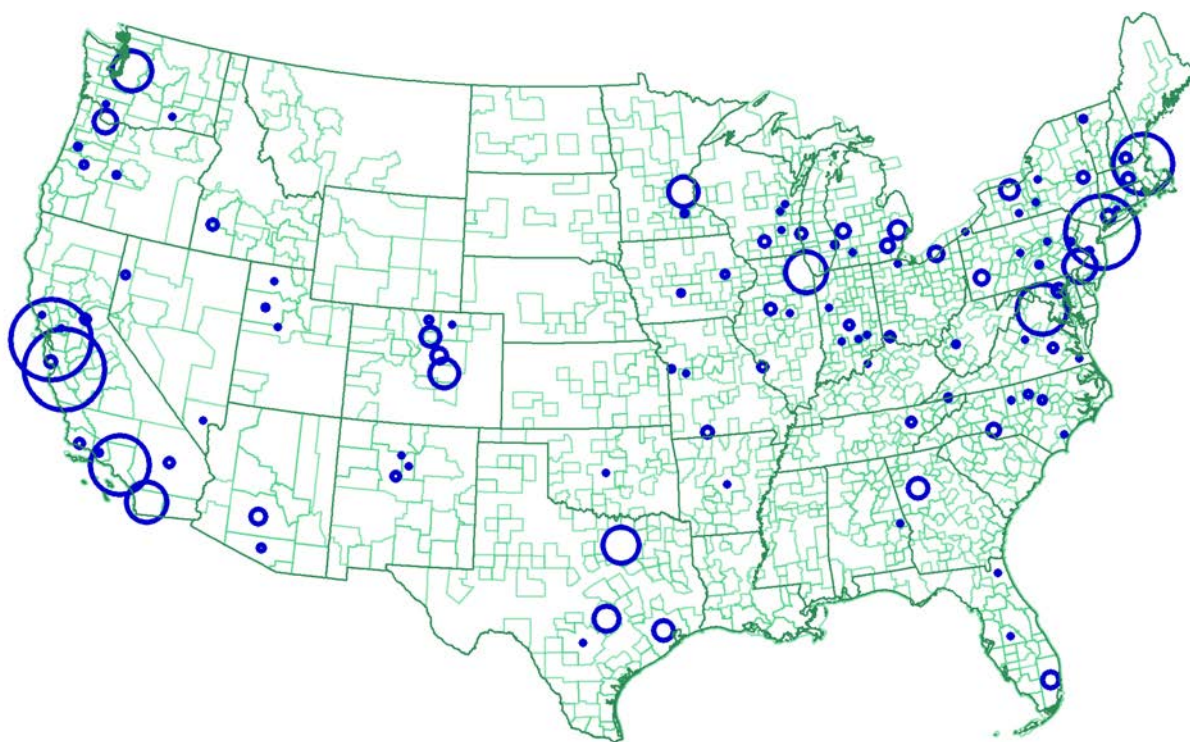
Figure 2 – Earnings calls and job postings for new technologies



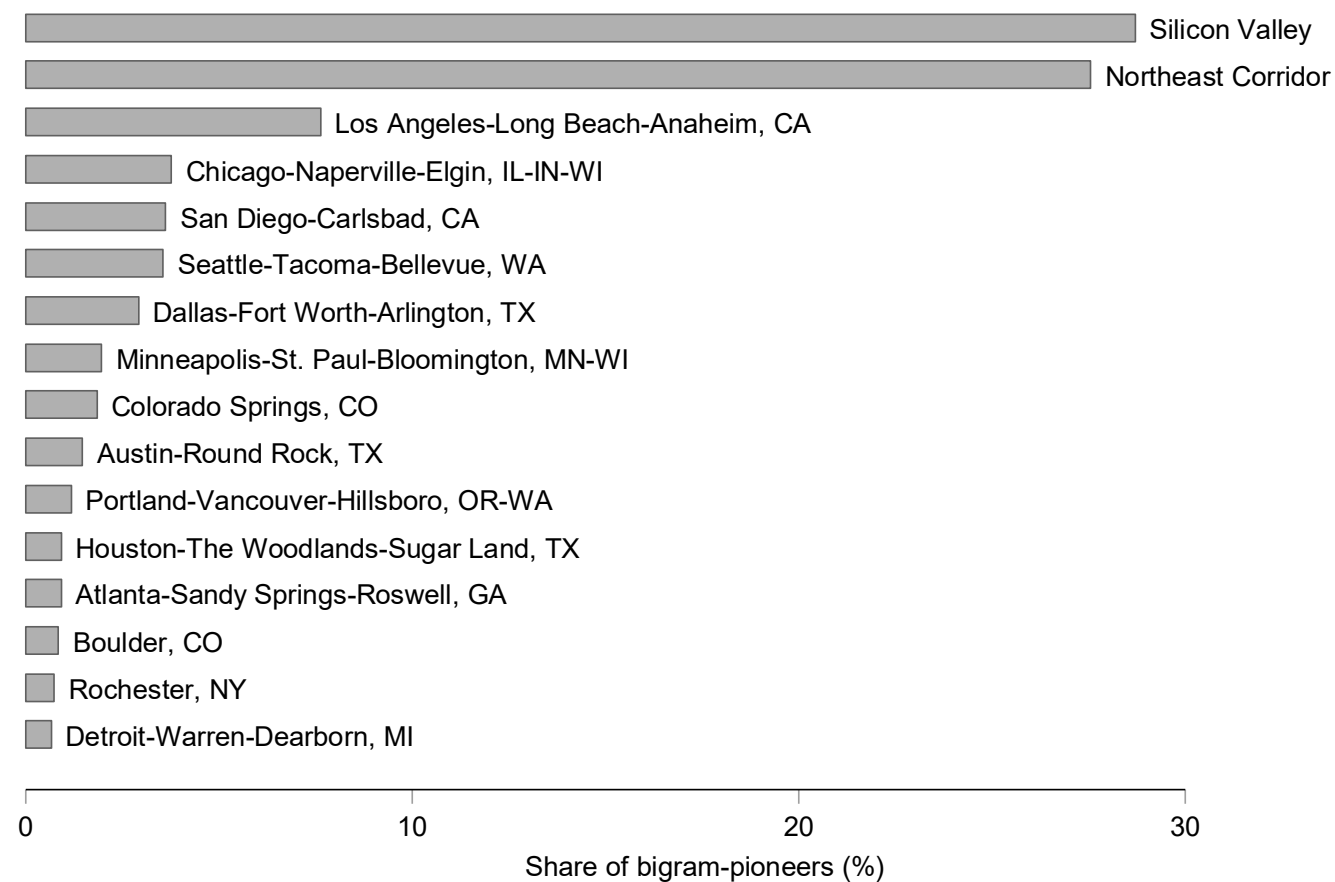
**Notes:** The figure shows a binned scatterplot at the technology bigram level of the number of the number of earnings calls that mention a given technology (y-axis) against the number of job postings that mention the technology bigram (x-axis). Some examples are labeled next to their bins.

**Figure 3 – Distribution of pioneer locations**

**Panel A: Pioneer Locations**



**Panel B: Distribution of Pioneer Locations**

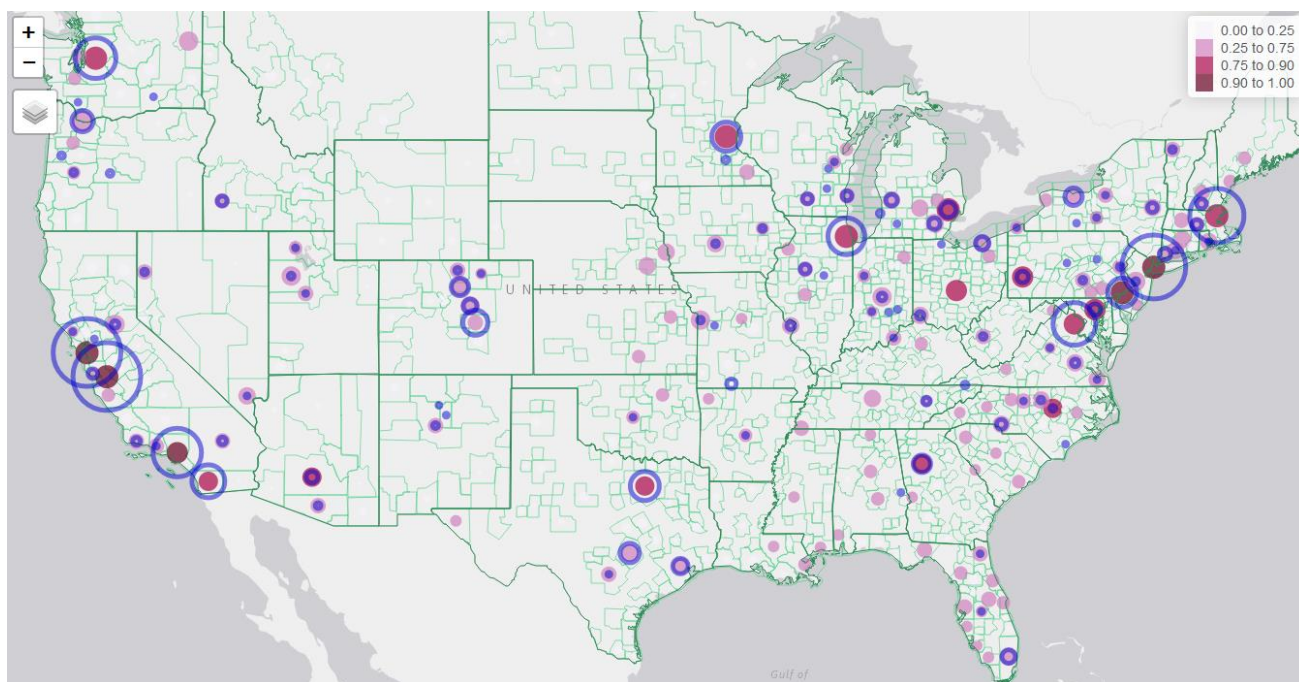


**Notes:** This figure shows the distribution of pioneer CBSAs. Panel A displays as blue circles CBSAs that are pioneer locations for at least one disruptive bigram. The size of the circles is proportional to the share of technology bigrams for which the CBSA is a pioneer location. Panel B shows a plot of the percentage of technology bigram-pioneer location pairs accounted for by each CBSA, for the top 16 CBSAs. We combine the CBSAs San Jose-Sunnyvale-Santa Clara, CA and San Francisco-Oakland-Hayward, CA, and label the region as Silicon Valley. Similarly, we combine New York-Newark-Jersey City, Boston-Cambridge-Newton, Washington-Arlington-Alexandria, and Philadelphia-Camden-Wilmington, and label the region as the Northeast Corridor.

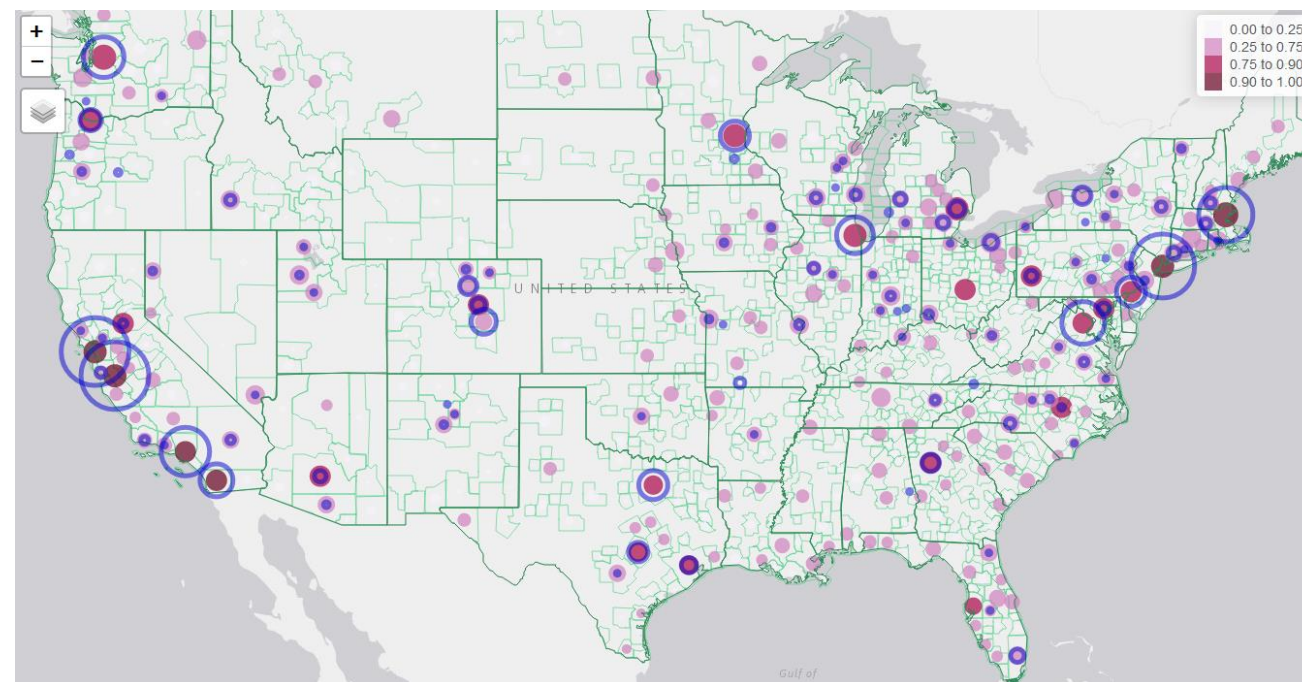


**Figure 4 – Geographic diffusion of technology job postings, by year since emergence**

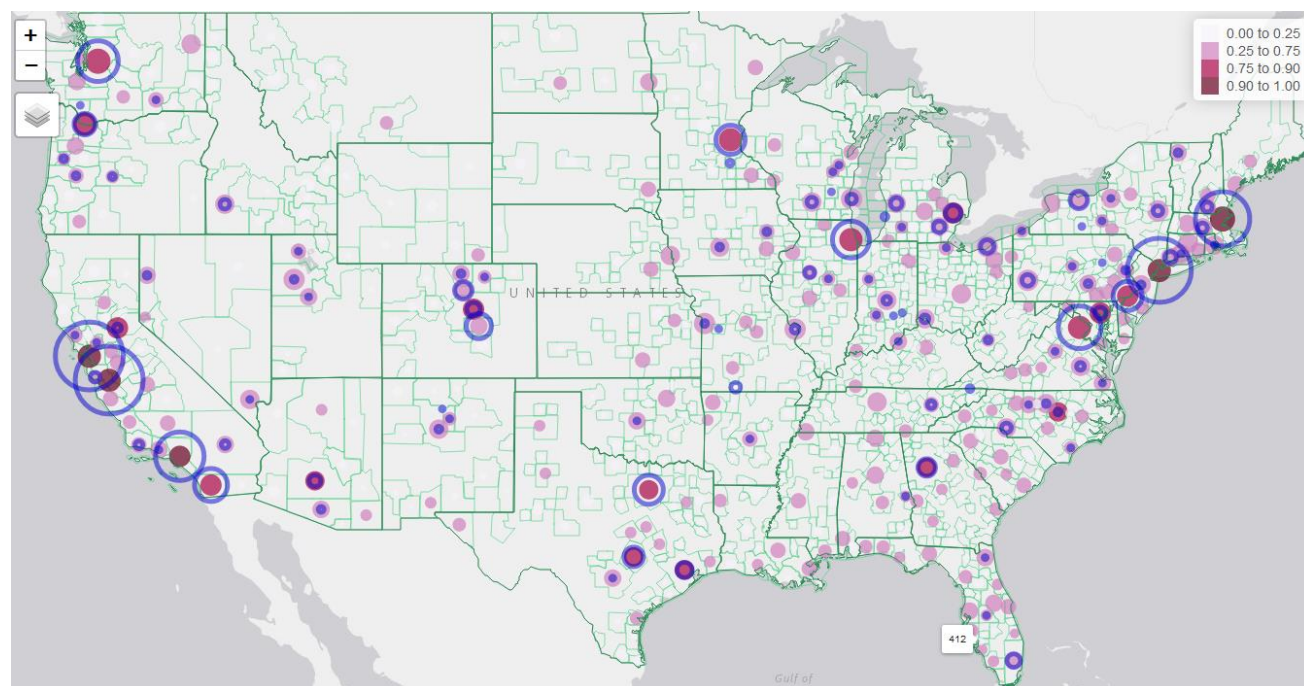
**Years since emergence: 0-5**



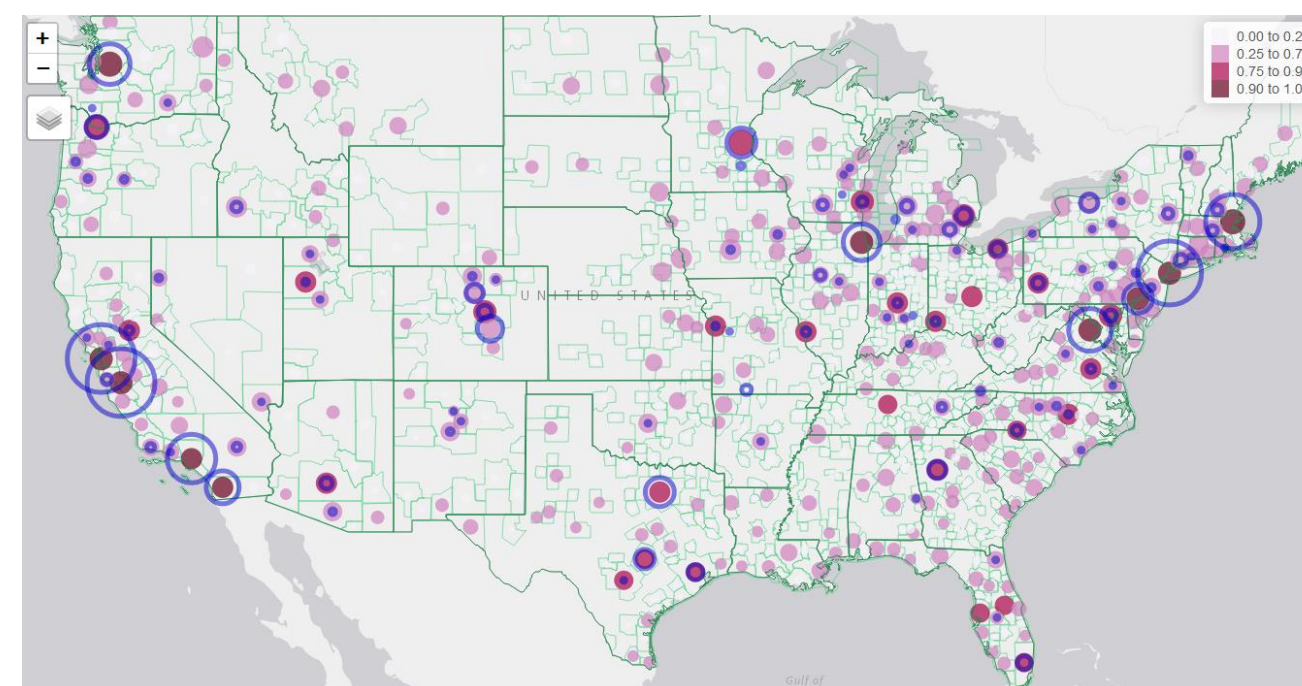
**Years since emergence : 6-10**



**Years since emergence: 11-20**

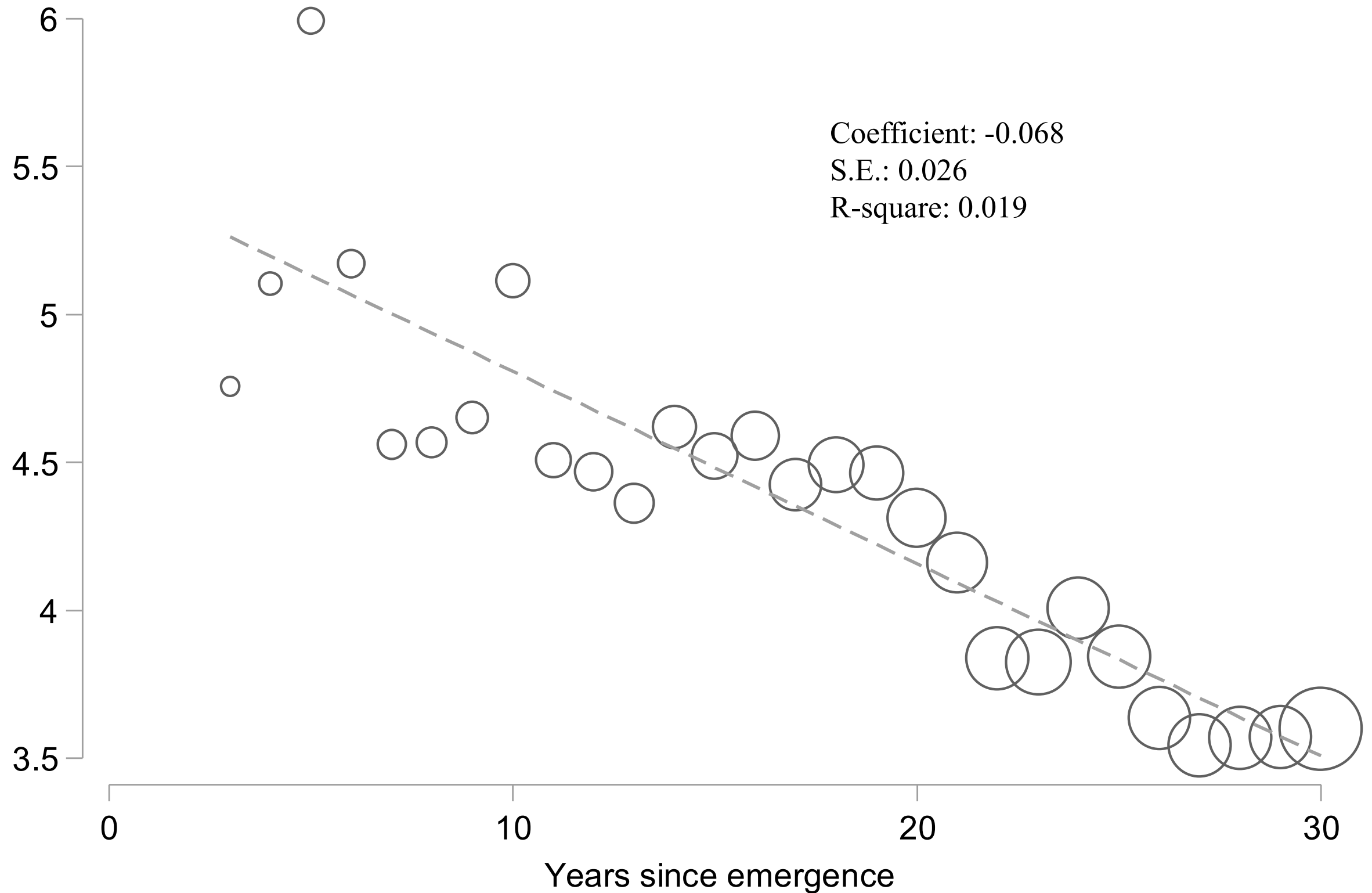


**Years since emergence : 21-30**



**Notes:** This figure plots maps with pioneer locations and job postings associated with technology bigrams by year since technology bigram emergence. Pioneer locations are marked with solid blue circles and technology job postings are in solid purple circles. For each CBSA and emergence year, we calculate the share of technology bigrams for which the CBSA records a non-negligible presence of technology jobs ( $Normalized\ share_{c,\tau,t} \geq 10\%$ ) and denote a higher share of technology bigrams with a darker color. For example, the first map plots the share of technologies with a normalized share of technology job postings greater than 10% for each CBSA between zero and five years since the emergence of the technology. The second map replicates this picture for six to ten years after emergence, and so forth. The sample for this map only contains technologies that appear in at least 100 earnings calls.

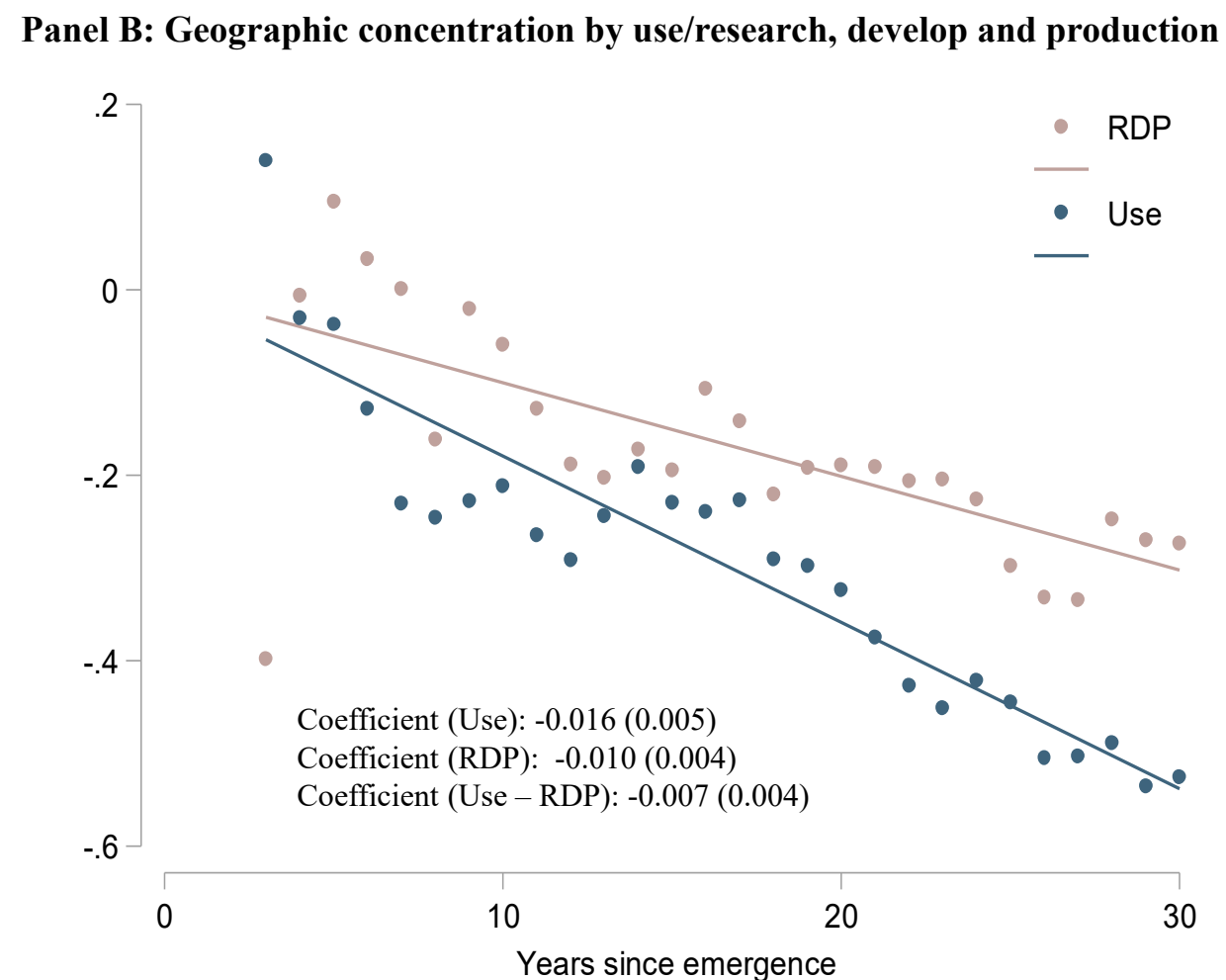
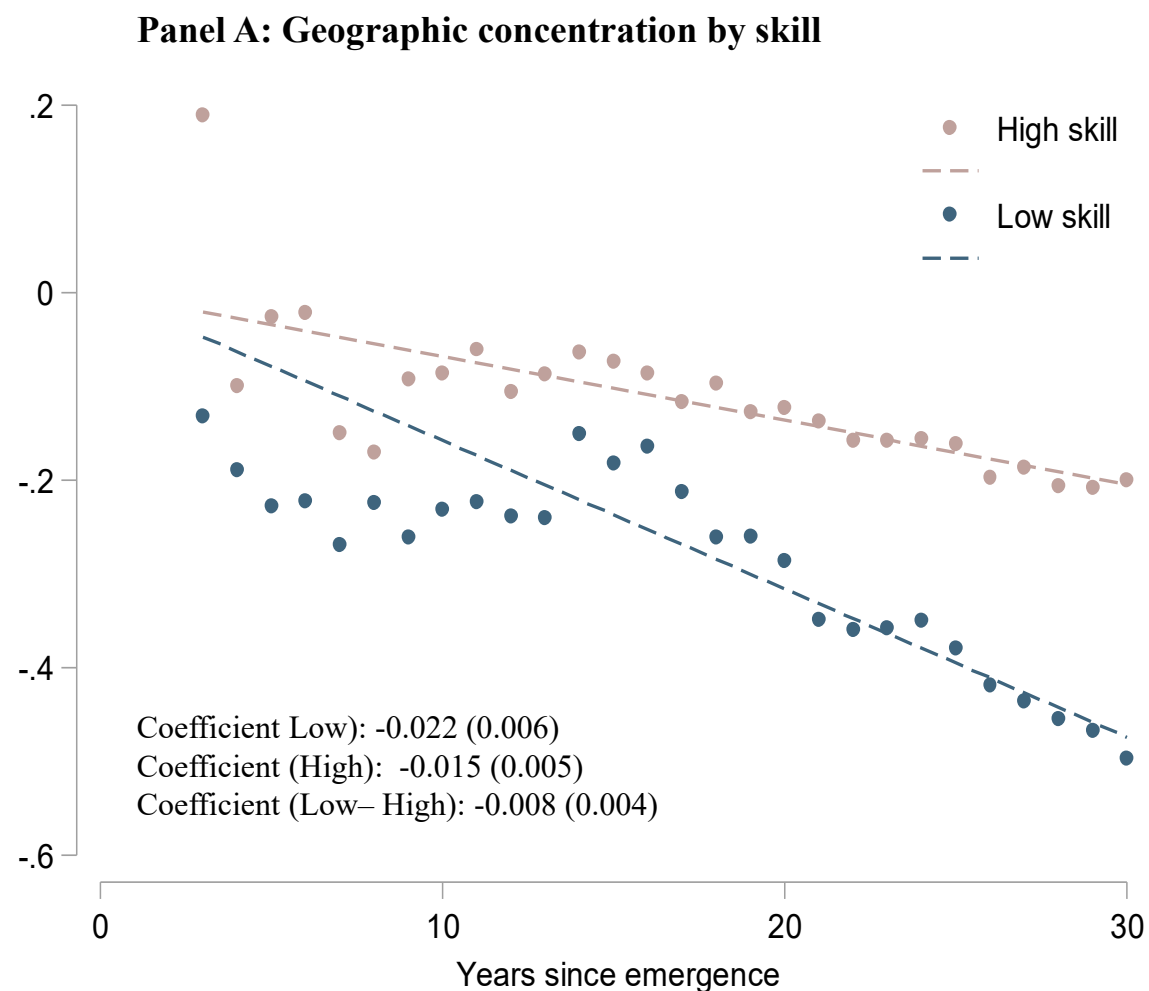
**Figure 5 – Geographic concentration of technology job postings across CBSAs, by year since emergence**



**Notes:** This figure shows a binned scatter plot at the technology bigram x year level of the coefficient of variation (CV) of the normalized share of technology job postings over time. We calculate the CV of the normalized share of technology job postings by dividing the standard deviation of  $Normalized\ Share_{c,t}$  across locations  $c$  in year  $t$  by its mean in year  $t$  for each technology bigram  $\tau$ . Each dot represents the weighted average of the CV (calculated across technologies) for each year since emergence, where the weight is the square root of the number of job postings for a bigram in a year, capped at 100. The circle sizes are proportional to the same weight. The regression line in the plot corresponds to a regression of the CV on year since emergence, as in Table 4, Panel A, column 1. We only include technology bigrams that appear in at least 100 earnings calls. Observations in and after the year of emergence are included.



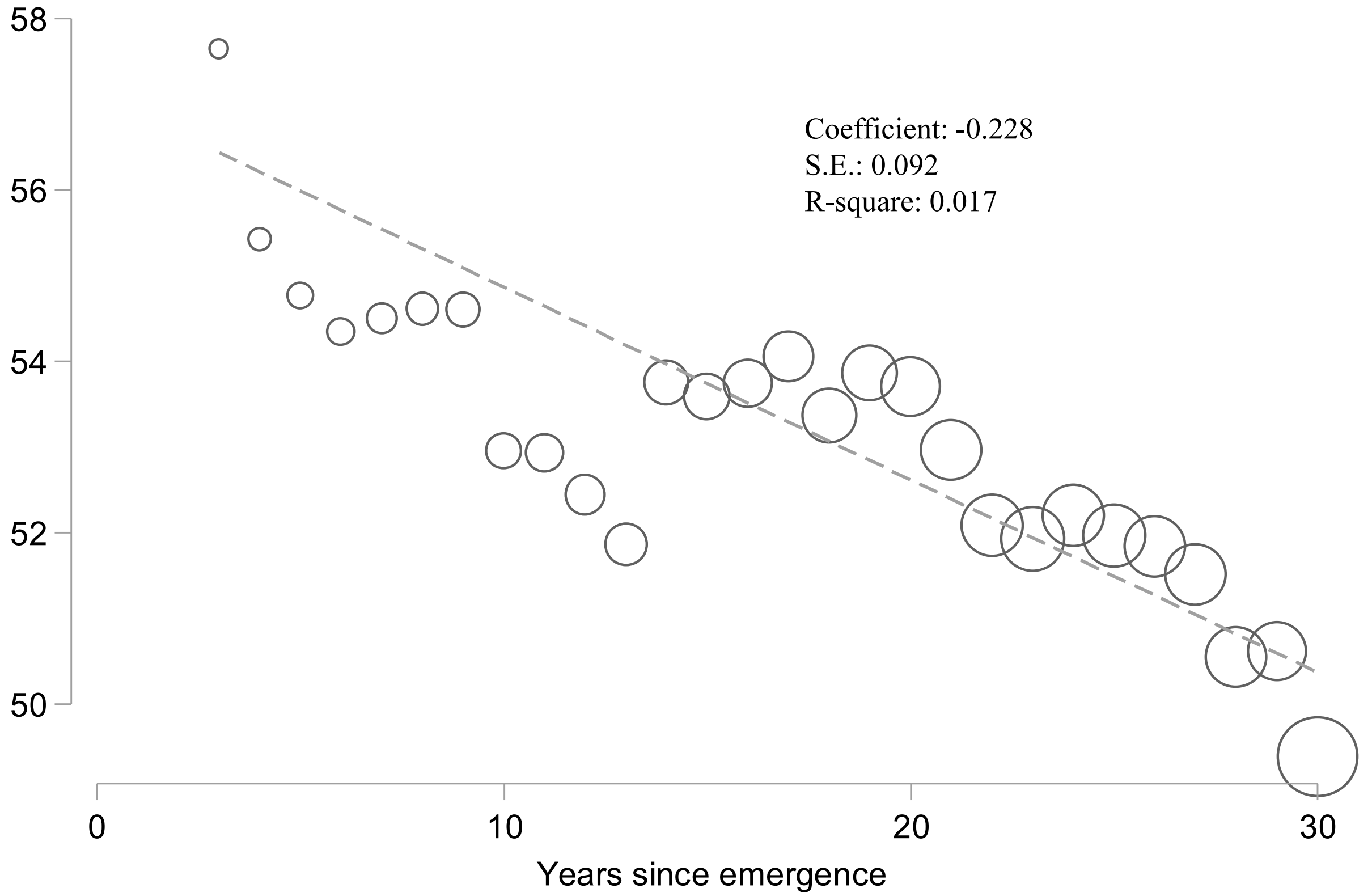
**Figure 6 – Geographic concentration relative to year since emergence, by skill and type of job posting**



**Notes:** This figure plots a binned scatter plot at the technology bigram x year x skill category level of the log of the coefficient of variation (CV) of the normalized share of bigram job postings relative to the year since emergence, by skill-level of job posting (high and low). The CV is calculated as the ratio of standard deviation to the mean of the normalized share across CBSAs for each technology bigram x year x skill category triplet. The red dots represent high-skill postings, and the blue dots represent low-skill postings. The fitted lines weigh observations by the square root of the total number of postings for a technology bigram in a year, capped at 100.

**Notes:** This figure plots a binned scatter plot at the technology bigram x year x job type level of the log of the coefficient of variation (CV) of the normalized share of bigram job postings relative to the year since emergence, by job type (RDP and use). The CV is calculated as the ratio of standard deviation to the mean of the normalized share across CBSAs for each technology bigram x year x job type triplet. The red dots represent RDP-related postings, and the blue dots represent use-related postings. The fitted lines weigh observations by the square root of the total number of postings for a technology bigram in a year, capped at 100.

**Figure 7 - Share of technology job postings requiring a college education, by year since emergence**



**Notes:** This figure shows a binned scatter plot at the technology bigram x year level of the share of technology postings requiring a college education by year since emergence. The share of college-educated postings for each technology bigram x year observation is measured as discussed in Section 5. Each dot represents a weighted average over technology bigrams of the share for each year since emergence, where the weight is the square root of the number of postings for a bigram in a year, capped at 100. The circle sizes are proportional to the same weight. The regression line in the plot corresponds to a regression of share of college-educated postings on the year since emergence as in Table 7, Panel A, column 1. We only include technology bigrams that appear in at least 100 earnings calls.

**Internet Appendix for**

**The Diffusion of New Technologies**

Aakash Kalyani, Nicholas Bloom, Tarek Hassan, Josh Lerner,  
Marcela Carvahlo, and Ahmed Tahoun

October 30, 2023

## **Data Appendix**

We process four sources of text data, and then combine them with U.S. government data to conduct our analyses.

### **1. Text Sources**

#### **1.1 Patents**

We download two separate sets of data for about six million utility patents, applied for at the U.S. Patent and Trademark Office (USPTO) between 1976 and 2014 and granted by 2018. First, we download full patent text XML files from the USPTO website. Second, we download processed patent variables, such as assignee names, inventor names and location, application and award year, citations (through 2018), and the primary CPC class, all from PatentsView.org. We restrict our sample to the sample of patents filed by U.S. inventors, about one-half of these awards.

We map the FIPS county identifier provided for inventors of each patent to Core-Based Statistical Areas (CBSAs), using a crosswalk provided by the U.S. Census Bureau. For patents with multiple inventor CBSAs, we assign the patent to each CBSA. We also standardize citation counts to control from truncation and time differences: we divide citations for each patent by the average number of citations for patents with a primary assignment to the same four-digit Combined Patent Classification (CPC) patent class and application year.

A typical patent award has six text sections: (1) title, (2) abstract, (3) background, (4) summary, (5) detailed description, and (6) claims. After we remove stop words (such as “of,” “the,” and “from”) following Kelly et. al. (2021) and Gentzkow et. al. (2019), we combine text from all of these sections into one large text string by appending all available text sections. We then break this large text string down into two-word combinations or bigrams. During this process, we use only those bigrams which are mentioned at least twice in a patent.

#### **1.2 Corpus of Historical American English**

The Corpus of Historical American English (COHA) is a collection of 116,759 documents published between 1880 and 1970. We download COHA from [www.english-corpora.org/coha](http://www.english-corpora.org/coha). These



include fiction and non-fiction books and newspaper and magazine articles. As with patents and earnings calls, we decompose these documents into about 400 million unique bigrams.

### **1.3 Wikipedia pages**

For each of our 36,563 novel and influential bigrams, we search on Wikipedia using the Wikipedia search functionality. We then download the full text of the first suggested Wikipedia page using the Wikipedia Python API (Wikipediaapi). For each of these Wikipedia pages, we store separately the title, the section headings, and the text of each section. While counting bigrams in Wikipedia, we count singular, plural, and unigram versions of these bigrams. For example, when counting `smart phone`, we include counts for “smart phone,” “smarts phone,” “smart phones,” “smartphone,” “smartphones,” “smartsphone,” and “smartsphones.” We follow the same approach when we count these bigrams in other corpora.

### **1.4 Earnings conference call transcripts**

From Refinitiv EIKON, we collect the complete set of 321,189 English-language transcripts of earnings conference calls held from 2002 through 2019. Out of these, we drop 5,552 transcripts because we could not reliably match them to a company name in Compustat. We obtain a total of 11,992 firms and 301,294 firm x quarter observations. We count our bigrams in the full text of these earnings calls and collect the number of earnings calls that each of our bigrams is mentioned in.

### **1.5 Burning Glass job postings**

From Burning Glass (BG), we obtain about 200 million job postings posted online in the U.S. Similar to patents, job postings data comes in two sets. The first contains the full text coded in XML files, while the second contains processed information about each job posting (such as occupation codes). We undertake minimal processing of job postings’ textual data: (1) removing non-letter sections of job postings; (2) removing the top 50 and bottom 50 words from each job posting, as mentioned in Section 2; and (3) as a consequence of step (2), excluding any job posting with less than 100 words. We then perform word counts over the remaining text. As before, while counting bigrams in Burning Glass, we count singular, plural and unigram versions of these bigrams. For example, when counting “smart phone,” we include counts for “smart phone,”

“smarts phone,” “smart phones,” “smartphone,” “smartphones,” “smartsphone,” and “smartsphones.” We follow the same approach when we count these bigrams in other corpora.

## **2. Merging Burning Glass Occupations with American Community Survey (2015)**

We obtain occupation and location demographic variables from the 2015 American Community Survey (ACS), downloaded on March 9, 2020. We examine respondents who are at least 25 years old, and report at least one year of schooling and a non-zero annual wage. We calculate the “share of college-educated people” in a particular occupation by dividing the number of people who report a particular occupation and have at least three years of college education by the total number of people who report the occupation. We calculate the “share of post-graduates” in a particular occupation by dividing the number of people who report a particular occupation and have at least a masters’ degree by the total number of people who report the occupation. We calculate the average wage in the occupation by taking an average over all annual incomes of people reporting a particular occupation. As for locations, we calculate skill levels using reported locations in the ACS and following the same methodology as for occupations. We also obtain population data for each CBSA from the ACS by performing a sample-weighted count of people who reported to live in a certain CBSA.

We merge the occupation level data from the ACS with occupation level aggregates in BG using six-digit SOC codes. Data on some six-digit SOC codes are reported in aggregated form in the ACS: for example, data on the occupational code 17-2021 (agricultural engineers) are reported as 17-20XX, along with class 17-2031 and others. In these cases, we map the six-digit SOC codes in BG to their aggregated values in the ACS.

As we do for occupations, we calculate share of college-educated people for CBSAs by dividing the number of people who report a particular CBSA as their residence in the ACS and have at least three years of college education by the total number of people in the CBSA. Other skill measures are calculated similarly.

## **3. Matching job postings to occupations, locations, industries and firms**

Burning Glass codes job postings into occupations (Standard Occupational Classification (SOC) codes), locations (counties), and industries (North American Industrial Classification Codes). BG also extracts an employer name with these job postings.

There are 836 occupations with six-digit SOC codes and 312 industries with NAICS Codes in the sample. We map counties in the BG data to CBSAs and use them as the unit for our geographical analysis. We do so by using the crosswalk made available by National Bureau of Economic Research (NBER). In this process, we lose about 2.8% of the job postings.

In order to assign firms to job postings, we use the employer strings provided by BG, which are available for 42.1% of job postings. Furthermore, these employer strings are not standardized or cleaned. For example, there are employer strings of the form “Tesla Motors Gigafactory,” “About Tesla,” and “Tesla Incorporated.” We generate firm identifiers from these raw employer strings using a modification of the process in Autor et. al. (2020):

- 1) We search the raw employer string on Bing.com and store the top five search result links. For instance, for the employer string “Tesla Incorporated,” we get <https://www.tesla.com>, <https://en.wikipedia.org/wiki/tesla-inc>, <https://www.britannica.com/topic/tesla-motors>, <https://www.bloomberg.com/quote/tsla/us>, <https://www.marketwatch.com/investing/stock/tsla>.
- 2) We group two employer strings under a single identifier if they share at least two out of top five links in common with each other.

Using this process, we group together 477,583 employer strings in BG into 329,158 unique firm identifiers.

We also match these employers to patent assignees using string matching. We implement the following modified “Term Frequency — Inverse Document Frequency” (tf-idf) algorithm. To do so, we:

- a. Decompose employer and assignee strings into 5 letter combinations. For example, “Alphabet” is broken into: “alpha,” “lphab,” “phabe,” and “habet.”
- b. Calculate a term frequency, which is the frequency of the five letter combination in the string. In our example of “Alphabet,” each combination uniquely appears in the strings.

We calculate an inverse document frequency (idf), which is inverse of the frequency with which the combination appears in all strings of assignees and BG employers.

- c. We combine the term frequency (tf) with an inverse document frequency (idf) to obtain a vector of combinations for each string:

$$v_{s,c} = tf_{c,s} * idf_c$$

where  $tf_{c,s}$  is the term frequency of the 5-letter combination  $c$  in string  $s$ ,  $idf_c$  is the inverse document frequency of each combination, and  $v_{c,s}$  is the value attributed to each combination separately for every string.

- d. Finally, we normalize each vector  $v_s$  so that the norm is 1. We then calculate similarities between two strings  $s$  and  $s'$  using dot product of their respective normalized vectors.

$$d_{s,s'} = v_s \cdot v_{s'}$$

- e. We match two strings if  $d_{s,s'} \geq 0.75$ .

A human audit of these matches resulted in an 86% accuracy rate.

#### 4. University Data

We download data on U.S. research universities from the U.S. National Science Foundation's Higher Education Expenditure on R&D (HERD) survey, which collects detailed statistics on research expenditure by these universities, and from the Integrated Postsecondary Education Data System (IPEDS) surveys provided by the U.S. Department of Education's National Center for Education Statistics (NCES). From these datasets, we construct the following variables:

- 1) Number of research universities in a CBSA: HERD provides details of universities which spend more than \$150,000 in research. We map university zip codes, provided as a part of university addresses, to CBSAs using a crosswalk provided by U.S. Census Bureau. Finally, we count the number of research universities in a CBSA to construct our variable.
- 2) University assets in a CBSA: IPEDS provides details of finances for most post-secondary educational institutions in the U.S. As with 1) above, we assign these universities to CBSAs and then aggregate their assets over CBSAs.

## Appendix Tables and Figures

Appendix Table 1 – Examples of novel influential bigrams from patents that pass and fail the Wikipedia technology bigram criteria

Technologies (pass)	Non-technologies (fail)
antigen binding	account manager
based solutions	active directory
branched chain	airway pressure
cell lines	aqueous phase
cloud computing	business model
communication channel	customer care
computer network	customer relationship
data quality	customer requests
data sources	dosage form
disk drive	email address
distributed computing	hardware software
fiber optic	healthcare provider
global positioning	homeland security
internet explorer	host computer
machine learning	identity theft
microsoft office	image data
mobile devices	management software
monoclonal antibody	mission critical
multiple access	performance management
optical fiber	performance metrics
personal computer	pharmaceutical composition
polymerase chain	pressure differential
positioning system	retail environment
programmable logic	risk management
scripting languages	search criteria
semiconductor devices	service offerings
temperature sensor	software engineering
user interface	transitory computer
vapor deposition	uninterruptible power
wireless technology	window assembly

**Notes:** This table lists examples of bigrams that pass (technologies) or fail (non-technologies) the Wikipedia filter. All examples are among the top 100 bigrams in terms of citation-weighted patent counts or number of job postings. Failure examples are sampled to reflect different sources of rejection in the Wikipedia filter.

Appendix Table 2 – Technology excerpts from earnings calls

Company	EC month	Excerpt
Ambarella Inc	4/2018	results that are many times higher in terms of processing performance per watt In March we successfully demonstrated to customer and investors our fully  AUTONOMOUS VEHICLE or embedded vehicle autonomy on Silicon Valley Road EVA navigated various traffic scenarios presented by Silicon Valleys challenging urban environment The fully autonomous
Cloudera Inc	4/2019	combined company road map which we rolled out in March of this year During this period of uncertainty we saw increased competition from the  PUBLIC CLOUD  vendors Second the announcement in March of Cloudera Data Platform our new hybrid and multicloud offering created significant excitement within our customer base CDP
NVIDIA Corp	7/2015	lot of very exciting development and were working with a lot of them because we have a platform that was really designed to fuse  COMPUTER VISION  cameras from all around the car as well as radars and LIDARS and sonars and be able to do path planning and all of
Proto Labs Inc	1/2015	orders in addition we added capacity to our manufacturing facility in europe in we completed our first acquisition purchasing fineline an  ADDITIVE MANUFACTURING  or 3D printing company based in raleigh north carolina the acquisition was completed last april and is highly complementary to proto labs roughly of our customers use
Collectar Biosciences Inc	10/2017	collaboration with Acunova Therapeutics each provide these types of strategic benefits Avicenna provides us with the unique opportunity to collaborate with experts in the antibody  DRUG CONJUGATE  or ADC field Not only does this provide the opportunity to work with a very promising small molecule payload but it also allows
L-3 Communications Holdings Inc	10/2002	metal detectors where they always make you take your shoes off This is a passive scanner as I told some of you It uses  MILLIMETER WAVE  It is nonintrusive and causes no harm or disease It will guarantee you won't have a weapon on you of any kind or be
InvenSense Inc	7/2016	as they strive to enable improved locationbased services and mapping user experience A significant opportunity for increasing our mobile content is UltraPrint our ultrasonic  FINGERPRINT SENSOR  I am very pleased to report that we are on track with the development of this gamechanging technology and have successfully passed several technology
SunPower Corp	10/2006	then be able to participate in the global electricity market which is measured in the form of trillion We have direct control over the solar cell and  SOLAR PANEL  portions of the value chain the technology core of the value chain that represents to of total installed costs in these
Donnelley Financial Solutions Inc	4/2018	speed and improve both the quality and consistency of business results for our clients in capital markets through the introduction of  MACHINE LEARNING  and artificial intelligence we will improve the efficiency of XBRL tagging and align with the efforts at the SEC to move from documents to data This investment

**Notes:** This tables presents examples of earning calls excerpts (in column 3) with 25 words before and after the mention of a technology bigram, with the firm (in column 1) and the date of the earnings call (in column 2).

Appendix Table 3 – List of Wikipedia titles associated with disruptive technology bigrams, by mentions in earnings calls

Mobile device	Bank account*	Music Player Daemon*	Genetic marker	Electronic discovery
Mobile phone	Solar panel	Laser diode	Combined Charging System	Linux powered device
Adverse event*	Communication channel	Nearfield communication	Vertical launching system	Drug design
User interface	Software defined radio	Power semiconductor device	Low Pin Count	Viral vector
Financial instrument*	Predictive modelling*	Thermography	Ion channel	Growth medium*
Smartphone	Semiconductor device	Heat treating	Plasma display	Sensor fusion
Wireless network	Single instruction multiple data	Performance appraisal*	Computer network	Pattern recognition
Active users*	Domain name*	Microsoft Access*	Electronic stability control	Input device*
Flat panel display	Barcode	Client computing	Wireless sensor network	Video capture
Digital content*	Project delivery method*	Dataspaces	Search algorithm*	Display device
Combination therapy	Distributed antenna system	Leased line*	Active noise control	Gesture recognition
Digital video	Sport utility vehicle	Dental implant	Urine collection device	Caregiver*
Model organism	Payment card	Multifunction printer*	Tunable laser	Policy entrepreneur*
Monoclonal antibody	Minimally invasive procedure	Tablet computer	Fingerprint*	Autonomous system Internet
Data source name*	Extremely high frequency	Semiconductor memory	Channel access method	Emergency Broadcast System*
Keypad	Central processing unit	DNA sequencing	Peritoneal dialysis	Oriented strand board
Social networking service	Location based service	Optical module	Error detection and correction	Variable computer science*
Digital camera	Lithium ion battery	Software versioning	Erotic target location error*	Robotic arm
Wireless LAN	Microsoft Windows*	User profile	Electronic game	Image organizer
Wireless	Unstructured data*	Data loss*	Computer simulation	Barcode reader
Transaction account	Compression seal fitting*	USB flash drive	Transdermal patch	Soulseek
Flash memory	LED circuit	Touchpad	File format*	Optical medium
Data communication	Baseband*	Powder coating	Semiconductor device fabrication*	Epidermal growth factor
Networking hardware	Digital image processing	Rapid prototyping	Dynamic random access memory	Cloud computing
Information privacy*	Video clip*	Continuous glucose monitor	Internet filter	Western blot
Online game	Multilayer perceptron*	Hematopoietic stem cell transplantation	Motion capture	Windowing system
Hard disk drive	Airborne early warning and control	Mental chronometry*	Solar thermal energy	Reduced relative clause*
Video game console	Cell site	Serverless computing*	Speed networking	Email
Mobile computing	Laptop	Liquid crystal display	Atomic layer deposition	
Digital television	Data model*	Silicon germanium	Fluid dynamics	
Media Player Windows*	Network virtualization	Data store	Drug eluting stent	
VLC media player*	File system*	Aggregate data*	Meta analysis*	
Virtual reality	Unmanned aerial vehicle	Ecommerce	Channel blocker	
Immortalized cell line	Portable media player	Data type*	Typical versus maximum performance*	
Personal computer	Multicore processor	Google Search*	Flexible electronics	
Nature based solutions*	Effective dose radiation*	Light emitting diode	SMS	
Selective Service System*	Web feed	Server computing	Video editing	
Network operating system	LED lamp	Component based software engineering	Digital signature	
Digital imaging	Facial recognition system	Software distribution	Electrooptical sensor	
Programming tool	Printing	Microcontroller	Identity verification service	
Data quality*	Antibody drug conjugate	Fiberoptic cable	Selective serotonin reuptake inhibitor	

Self driving car	Insulin pump	List of search engines	CMOS
Passive optical network	Voicemail	Push technology	Optical disc drive
Augmented reality	Speech recognition	Machine learning	Magnetic resonance imaging
Smart card	Satellite navigation device	Data file*	Amazon Relational Database Service
Card security code	Building automation	Thermometer	Programming language*
Power electronics	User activity monitoring	Software architecture	Computer data storage
Text messaging	Data integrity	Dendritic cell*	Medical imaging
Home automation	Programmable logic device	Therapeutic drug monitoring	Ventricular assist device
Receptor antagonist	Web query	Genetic engineering	Dry powder inhaler
Instant messaging	Flow cytometry	Da Vinci Surgical System*	Network switch
Multimedia Messaging Service	Virtual private network	Mesh networking	Urine test strip
Coiled tubing	Desktop computer	Epidermal growth factor receptor	Label printer
Digital subscriber line	Network monitoring	Monoclonal antibody therapy	Distributed computing
Systems architecture*	Audio file format*	Communication protocol*	Physical media
OS level virtualization*	Lithium battery	Neural network	Glucose meter
Software Updater	Global Positioning System	Docking station	Electronic program guide
Patch computing	Visitor Based Network	Recombinant DNA	Adeno associated virus
Expected value*	Active site*	Duty cycle*	List of copy protection schemes*
Agent based model	Memory card	Mass spectrometry	Debit card
Data access object*	Gallium nitride	Data compression*	Laser scanning
Hybrid electric vehicle	Load balancing computing	Cruise control	Intensive care unit*

**Notes:** This table lists all technologies (Wikipedia titles) associated with technology bigrams which are mentioned in at least 100 earnings calls. \* next to a technology indicates that it is dropped in human audit described in Section 7.



Appendix Table 4 – Top technologies by pioneer states, by share of early patenting

Region	Division	State Name	Title	Pct. Early Patents	Year of Emergence
(1)	(2)	(3)	(4)	(5)	(6)
Northeast	New England	Connecticut	Label printer	6.19	1986
		Massachusetts	Antibody-drug conjugate	13.58	1998
		New Hampshire	Computer data storage	14.3	1992
	Mid Atlantic	New Jersey	Transaction account	19.45	2000
		New York	Digital imaging	16.95	1992
	Pennsylvania	Machine learning	17.12	1992	
Midwest	East North	Illinois	Ventricular assist device	12.59	1983
		Indiana	Hybrid electric vehicle	8.61	1995
		Michigan	Electronic stability control	49.93	1996
		Ohio	Da Vinci Surgical System	17.5	2005
	Wisconsin	Peritoneal dialysis	15.55	1990	
	West North	Minnesota	Optical disc drive	14.15	1991
Missouri		Sensor fusion	24.64	1996	
South	South Atlantic	Florida	Dental implant	12.87	1987
		Georgia	Oriented strand board	8.34	1991
		Maryland	Adeno-associated virus	13.28	1994
		North Carolina	Gallium nitride	24.36	1993
	Virginia	Selective Service System	7.25	1990	
	East South	Tennessee	Lithium battery	4.79	1986
	West South	Texas	Coiled tubing	51.74	1988
West	Mountain	Arizona	Cruise control	4.71	1993
		Colorado	Network monitoring	23.1	1991
		Idaho	Dynamic random-access memory	10.11	1991
		Utah	Fluid dynamics	5.62	1993
		Pacific	California	Continuous glucose monitor	73.55
Oregon	Multifunction printer		29.59	1997	
		Washington	Distributed computing	25.99	1987

**Notes:** This table reports top technology (in column 4) by share of early patenting for a sample of states (in column 3). We report the corresponding census division in column 2, and the census region in column 1. The sample of states are those which file at least 50,000 patents during our sample period. In column 6, we report the share of early patenting, calculated as  $Share\ Early\ Patenting_{\tau,s} = \frac{Early\ Patents_{\tau,s}}{\sum_j Early\ Patents_{\tau,j}}$  for every technology  $\tau$  and state  $s$ . We drop the titles “policy entrepreneur,” “caregiver,” and “reduced relative clause” that pass the Wikipedia criteria of technology bigrams but are removed in the human audit of technology bigrams.

Appendix Table 5 – Top occupations for technology - “Machine Learning”

Occupations	Pct. “Machine Learning” postings	Total Job Postings
Computer and Information Research Scientists	40.83	179,636
Astronomers	3.33	10,236
Life Scientists, All Other	2.41	27,691
Computer Hardware Engineers	2.38	91,259
Statisticians	2.31	193,339
Computer Science Teachers, Postsecondary	2.24	34,221
Operations Research Analysts	1.80	874,226
Database Administrators	1.62	1,110,414
Social Science Research Assistants	1.56	53,040
Biological Scientists, All Other	1.55	97,722
Software Developers, Applications	1.50	6,963,792
Physical Scientists, All Other	1.21	12,176
Engineering Teachers, Postsecondary	1.17	13,606
Social Scientists and Related Workers, All Other	1.15	59,306
Detectives and Criminal Investigators	1.11	130,801
Biomedical Engineers	1.06	17,543
Computer and Information Systems Managers	1.04	205,267
Financial Specialists, All Other	1.03	290,386
Architectural and Engineering Managers	1.01	586,301
Computer Occupations, All Other	1.00	6,172,457

**Notes:** This table lists the top occupations (in column 1) by the share of postings that are associated with technology “Machine Learning” (in percent, in column 2). Column 3 reports the total job postings for the occupation. The technology “Machine Learning” refers to Wikipedia title “Machine learning” and the associated bigrams “machine learning” and “learning algorithms.”

Appendix Table 6 – Human audit results of technology job postings

Panel A: Audit results			
Audit	Use	Produce	Total
Describes company	6%	10%	16%
Describes task	46%	34%	80%
Neither	NA	NA	4%
Panel B: Audit results after clipping top 50 and bottom 50 words			
Audit	Use	Produce	Total
Describes company	2%	2%	4%
Describes task	55%	36%	91%
Neither	NA	NA	5%
Panel C: Examples excerpts			
Produce	“we are looking for a ux developer to join our fast growing team our mission is to make your interactions with <b>touchscreens</b> more interesting more natural and more engaging we have developed a novel haptic <b>touchscreen</b> technology that not only tracks the fingertips but controls what they feel”		
Use	“duties of this job the employee is required to occasionally use clarity of vision at approximately feet or more stoop bend the body downward and forward by the spine at the waist frequently use a keyboard, key <b>touch screen</b> , or mouse to enter text or data into a computer”		
Neither	“our super cool office space which doesn’t feel like an office is designed with our employees in mind techy surroundings a great outdoor space with Wi-Fi hookups for your laptop plus Bluetooth capabilities for <b>music streaming</b> we enjoy cultivating a supportive and all around positive culture that keeps our employees happy this will be a place you will want to come to everyday”		

**Notes:** This table presents the results from a human audit of Burning Glass technology job postings. As a part of the human audit, we classify each of 1,000 randomly sampled job postings into two ways: 1) whether the technology reference in the job posting refers to the company or the task content of the job posting, and 2) whether the job describes the use or the production of the technology. See the main text for details. In Panel A, we perform the audit on the original text of job postings for a smaller set of 100 postings. In Panel B, we clip the text of job postings by 50 words at the top and bottom, resample 1,000 postings, and then repeat the audit. Panel C provides examples of job postings classified under the “use,” “produce,” and “neither” categories.

Appendix Table 7 – Summary statistics

Variable	(1) N	(2) Mean	(3) SD	(4) p25	(5) p50	(6) p75
<b>Panel A: Technology bigrams</b>						
Postings	1,899	26633.79	127987.98	45.00	534.00	6440.00
Earnings Calls	1,899	132.86	471.41	1.00	13.00	85.00
Cite wt. Patents	1,899	4643.53	7289.85	1506.59	2513.45	5140.56
Emergence Year	1,899	1991.54	5.28	1989	1991	1994
# Pioneers	1,899	5.90	2.60	4	6	7
<b>Panel B: Technology bigram by year</b>						
Tech Postings	15,498	3370.21	15528.70	18	111	920
Coefficient of Variation	15,498	5.70	4.67	2.47	4.10	7.44
Share College Educated	15,512	54.82	13.06	47.29	56.46	63.26
Wage	15,512	62062.27	10910.91	55217.73	63263.68	69452.33
Share Research Postings	15,512	1.75	3.71	0.00	0.00	2.01
Share Produce Postings	15,512	5.36	9.55	0.00	0.00	7.74
Share RDP Postings	15,512	5.81	10.12	0.00	0.00	8.48
Share Knowledge Postings	15,512	7.50	11.84	0.00	0.00	12.23
<b>Panel C: CBSA by technology bigram by year</b>						
Total Postings	17,634,114	20670.29	81760.65	1252	2825	8897
Tech Postings	17,634,114	2.96	67.09	0.00	0.00	0.00
Normalized Share	14,391,504	0.58	2.77	0.00	0.00	0.28

**Notes:** This table shows summary statistics – the number of observations, mean, standard deviation, 25<sup>th</sup> percentile, median and 75<sup>th</sup> percentile - for variables used in the analyses of the paper. Panel A reports statistics for our sample of technology bigrams. Panel B reports summary statistics for our datasets at the technology bigram x year level, which are used in region broadening (Table 4) and skill broadening (Table 7) regressions. Panel C reports summary statistics for our CBSA x technology bigram x year data used in pioneer advantage regressions in Table 5. The normalized share of bigram jobs in all panels is calculated as  $Normalized\ share_{c,\tau,t} = \frac{share\ jobs\ exposed_{c,\tau,t}}{share\ jobs\ exposed_{\tau,t}}$ , where  $c$  is a location (CBSA),  $\tau$  is a technology, and  $t$  is calendar year. The coefficient of variation is calculated using the normalized share of technology job postings by dividing the standard deviation of  $Normalized\ share_{c,\tau,t}$  across locations  $c$  in year  $t$  by its mean in year  $t$  for each technology bigram  $\tau$ . Skill level variables (in Panel B – Share College Educated and Wage) are calculated using  $Skill_t^\tau = \frac{\sum_o N_{o,t}^\tau \chi_{o,2015}}{\sum_o N_{o,t}^\tau}$ , where  $N_{o,t}^\tau$  is the number of Burning Glass job postings mentioning technology bigram  $\tau$  that are in SOC code  $o$  at time  $t$ , and  $\chi_{o,2015}$  is the average skill level for occupation  $o$ , as measured by the 2015 ACS. The share of research/produce/RDP/knowledge postings are calculated for each technology bigram  $\tau$  and year  $t$  according to the process explained in Section 5a.

Appendix Table 8 – Robustness: Testing for non-linearity in ‘region broadening’

Panel A: Comparing linear, quadratic and logarithmic specifications				
	<i>Coefficient of Variation</i> <sub><math>\tau,t</math></sub>			$\log(\text{Coefficient of Variation}_{\tau,t})$
	(1)	(2)	(3)	(4)
<i>Years since emergence</i> <sub><math>\tau,t</math></sub>	-0.068*** (0.026)	-0.028 (0.133)	-0.305*** (0.076)	-0.015** (0.006)
<i>(Years since emergence</i> <sub><math>\tau,t</math></sub> ) <sup>2</sup>		-0.001 (0.003)	0.004** (0.002)	
Constant	5.577*** (0.645)	5.203*** (1.444)	9.008*** (0.861)	1.523*** (0.158)
R-squared	0.019	0.019	0.817	0.017
N	4,270	4,270	4,270	4,270
Bigram FE	NO	NO	YES	NO

Panel B: Testing for general non-linearity		
Specification: $\text{Coefficient of Variation}_{\tau,t} = \beta_0 + \beta_1 * (\text{years since emg}_{\tau,t})^{\beta_2}$		
Initial conditions: $\beta_0 = \text{mean}(\text{Coefficient of Variation}_{\tau,t}), \beta_2 = 1$		
$\beta_0$	$\beta_1$	$\beta_2$
5.134*** (0.375)	-0.009 (0.020)	1.536** (0.626)
H0: $\beta_2 = 1$	0.73	
p-value	0.392	
H0: $\beta_1 = \beta_2 = 0$	1,895.00	
p-value	0.000	

**Notes:** This table reports results from tests for non-linearity of the relationship between the *Coefficient of Variation* <sub>$\tau,t$</sub>  and *Years since Emergence* <sub>$\tau,t$</sub> . Panel A, column 1 reports results from a regression of the coefficient of variation on the year since emergence, the same specification as in Table 4, Panel A, column 1. Column 2 adds a quadratic term to this regression and column 3 adds technology bigram fixed effects. Column 4 presents results for a regression where we replace *Coefficient of Variation* <sub>$\tau,t$</sub>  with  $\log(\text{Coefficient of Variation}_{\tau,t})$  as the dependent variable. Panel B tests for non-linearity by estimating the specification:  $\text{Coefficient of Variation}_{\tau,t} = \beta_0 + \beta_1 * (\text{years since emg}_{\tau,t})^{\beta_2}$  with non-linear least squares and initial conditions  $\beta_0 = \text{mean}(\text{Coefficient of Variation}_{\tau,t})$  and  $\beta_2 = 1$ . The regressions are restricted to the sample of technology bigrams that appear in at least 100 earnings calls. For more details on the specification, refer to the note in Table 4.

Appendix Table 9 – Robustness: Pioneer location advantage in technology hiring

	<i>Normalized Share<sub>c,t,τ</sub></i>			
	(1)	(2)	(3)	(4)
<i>Rate of decline per year</i>	-0.027***	--0.027***	--0.024***	--0.024***
	0.003	0.003	0.005	0.005
<i>Implied years to zero advantage</i>	37.04	37.04	41.67	41.67
<i>Pioneer<sub>c,τ</sub></i>	1.321***	1.323***	1.118***	1.121***
	(0.254)	(0.255)	(0.294)	(0.295)
<i>Pioneer<sub>c,τ</sub> * Years since emg<sub>τ,t</sub></i>	-0.035***	-0.035***	-0.027**	-0.027**
	(0.011)	(0.011)	(0.012)	(0.012)
R-squared	0.030	0.037	0.033	0.040
N	7,751,024	7,751,024	7,751,024	7,751,024
Bigrams	835	835	835	835
Bigram FE	YES	NA	YES	NA
CBSA FE	YES	YES	NA	NA
Year FE	YES	NA	YES	NA
Bigram x Year FE	NO	YES	NO	YES
CBSA x Year FE	NO	NO	YES	YES
Std. Errors (cluster)	Wiki Title	Wiki Title	Wiki Title	Wiki Title

**Notes:** This table reports robustness checks for results from the pioneer advantage regressions in Table 5. We are focusing on the change in the rate of decline per year, so place these findings above the regression coefficients. Column 1 repeats the regression specification from Table 5, column 3. Columns 2, 3 and 4 add CBSA x Bigram fixed effects and CBSA x Year fixed effects iteratively. The first row reports the decay rate, as defined in Table 5. The regressions use all technology bigrams that appear in at least 1000 job postings in our sample. For more details on the specification, refer to the note in Table 5.

Appendix Table 10 – Top occupations for example technologies

Emergence Year	Technology	Top Occupations (by pct. share of postings)
1979	Hard disk drive	Broadcast Technicians (4.5); Computer Operators (2.7)
1980	Barcode reader	Packers and Packagers, Hand (0.5); Library Assistants, Clerical (0.2)
1981	Laser diode	Chiropractors (0.4); Physicists (0.3)
1982	Personal computer	Automotive Glass Installers and Repairers (26.7); Demonstrators and Product Promoters (22.6)
1983	Flatpanel display	Upholsterers (2.6); Audio and Video Equipment Technicians (2.0)
1984	User interface	Multimedia Artists and Animators (26.7); Web Developers (23.4)
1985	Mobile phone	Graders and Sorters, Agricultural Products (20.7); Advertising Sales Agents (15.8)
1986	Facial recognition system	Marriage and Family Therapists (1.8); Child, Family, and School Social Workers (1.6)
1987	Digital video	Audio and Video Equipment Technicians (7.4); Film and Video Editors (5.1)
1988	Model organism	Astronomers (2.3); Life Scientists, All Other (1.8)
1989	Mobile device	Electronic Equipment Installers and Repairers, Motor Vehicles (12.6); Foresters (11.3)
1990	Debit card	Pharmacy Aides (10.3); Railroad Conductors and Yardmasters (5.5)
1991	Flash memory	Computer Hardware Engineers (0.9); Radio, Cellular, and Tower Equipment Installers and Repairs (0.5)
1992	Machine learning	Computer and Information Research Scientists (56.7); Astronomers (4.0)
1993	Financial instrument	Financial Specialists, All Other (0.7); Economists (0.4)
1994	Active users	Advertising Sales Agents (2.0); Sales Representatives, Services, All Other (0.3)
1995	Hybrid electric vehicle	Electronics Engineers, Except Computer (0.2); Electrical Engineers (0.2)
1996	Digital content	Producers and Directors (3.6); Multimedia Artists and Animators (3.6)
1997	Multicore processor	Computer Hardware Engineers (0.2); Electronics Engineers, Except Computer (0.1)
1998	Information privacy	Information Security Analysts (2.7); Financial Specialists, All Other (0.7)
1999	Unmanned aerial vehicle	Avionics Technicians (2.5); Commercial Pilots (1.4)
2000	Transaction account	Traffic Technicians (0.2); Brokerage Clerks (0.2)
2001	Smartphone	Automotive Glass Installers and Repairers (28.5); Home Appliance Repairers (21.5)
2002	Online game	Fine Artists, Including Painters, Sculptors, and Illustrators (1.2); Multimedia Artists and Animators (0.8)
2003	Social networking service	Reporters and Correspondents (5.4); Radio and Television Announcers (3.4)
2004	Electronic discovery	Graders and Sorters, Agricultural Products (13.2); Parts Salespersons (1.3)
2005	LED circuit	Electrical and Electronics Drafters (0.2); Electronics Engineers, Except Computer (0.1)
2006	Augmented reality	Computer Hardware Engineers (0.3); Interior Designers (0.3)
2007	Self-driving car	Computer Hardware Engineers (0.5); Armored Assault Vehicle Crew Members (0.2)

**Notes:** This table lists the top two occupations and the percentage of postings (in column 3) that mention technology bigrams associated with top technologies (in column 2). The list of technologies is the same as our list in Table 1.

Appendix Table 11: Robustness: Region broadening by skill

	<i>Coefficient of Variation</i> $_{\tau,t}$		
	EC $\geq$ 100		
	(1)	(2)	(3)
Skill Level	Low	Medium	High
<i>Years since emergence</i> $_{\tau,t}$	-0.343*** (0.035)	-0.256*** (0.029)	-0.194*** (0.023)
Constant	16.715*** (0.795)	13.685*** (0.655)	12.207*** (0.512)
R-squared	0.791	0.784	0.778
N	4,188	4,188	4,188
Bigrams	427	427	427
Bigram FE	YES	YES	YES
Std. Errors (cluster)	Wiki Title	Wiki Title	Wiki Title
<b><i>Rate of decline per year (%)</i></b>	<b>-2.05</b>	<b>-1.87</b>	<b>-1.59</b>

**Notes:** This table reports region-broadening regressions for low-, medium-, and high-skill job postings at the technology bigram x year level. Columns 1, 2 and 3 show separate regressions for high-skill (column 1), medium-skill (column 2), and high-skill (column 3) job postings. For definitions of these concepts see Section 4.c of the main text. The regressions are restricted to the sample of technology bigrams that appear in at least 100 earnings calls. For more details on the specification, refer to the note in Table 4.



Appendix Table 12 – Keywords for research, development, and production (RDP) job postings

Keyword	Pct, Tech. Job Postings	Keyword	Pct, Tech. Job Postings
development*	15.30	the leading	0.18
engineering*	4.67	in creating	0.17
design and*	4.40	the worlds	0.12
develop and*	1.39	to market	0.10
to create	1.02	adoption of	0.10
and develop*	1.00	customization of*	0.09
range of	0.80	to customize*	0.03
leader in	0.79	enhancements of	0.02
designing and*	0.70	and experimentation	0.01
research and*	0.67	and exploit	0.01
developing and*	0.63	exploitation of	0.01
build and	0.59	in expanding	0.01
team of	0.57	and commercializing*	0.00
in developing*	0.54	to personalize*	0.00
and developing*	0.50	customized product*	0.00
a global	0.48	researching of*	0.00
high performance	0.39	and redevelopment	0.00
in building*	0.29		

**Notes:** This table reports the list of keywords used to classify technology job postings into those requiring research, development, and production (RDP) of the technology. We flag job postings that mention these RDP keywords within 15 words of a technology bigram as requiring the RDP of the technology. The sample used to calculate these shares only includes technology bigrams that appear in at least 100 earnings calls. The asterisk \* identifies keywords for research and development.

Appendix Table 13 – Robustness: Skill broadening robustness

	Share College Educated			
	Baseline (2010-2019)	W/ 2007	2010-2015	2016-2019
	(1)	(2)	(3)	(4)
Year since emergence	-0.288*** (0.079)	-0.294*** (0.072)	-0.260*** (0.088)	-0.228** (0.114)
Constant	59.095*** (1.794)	59.313*** (1.605)	58.901*** (1.814)	57.047*** (2.933)
R-squared	0.019	0.021	0.016	0.007
N	8,347	9,277	5,008	3,339
Bigram FE	NO	NO	NO	NO

**Notes:** This table reports robustness of skill broadening results. All regressions are at the technology bigram x year level. The dependent variable is the average share of job postings mentioning technology bigram  $\tau$  in year  $t$  that require a college degree. Column 1 uses our baseline sample of Burning Glass job postings from 2010-2019 (replicating Table 7, Panel A, column 2); column 2 adds 2007 to our baseline sample. Columns 3 and 4 split the sample into two: column 3 only includes 2010-2015 and column 4 only includes 2016-2019. The regressions use all technology bigrams that appear in at least 1000 job postings in our sample. For more details on the specification, refer to the note in Table 7.

Appendix Table 14 – Keywords for “training” job postings

Keyword (1)	Pct, Tech. Job Postings (2)	Keyword (1)	Pct, Tech. Job Postings (2)
knowledge of	16.888	proficiency of	0.045
experiences with	13.489	expertise of	0.034
understanding of	6.098	experience using	0.026
working knowledge	4.173	specialization in	0.025
responsible for	3.777	competence of	0.007
proficiency in	3.680	indepth knowledge	0.002
proficient in	3.626	credentials in	0.001
familiar with	3.190	knowhow of	0.000
skills in	1.432	critical of	0.000
experience with	1.398		
expertise in	1.236		
familiarity with	1.048		
knowledge in	0.876		
experience of	0.400		
extensive knowledge	0.388		
knowledge with	0.222		
grasp of	0.078		
understanding in	0.056		

**Notes:** This table reports (in column 1) the list of keywords used to classify technology job postings into those that require training in the technology. We flag job postings that mention these training keywords within 15 words of a technology bigram as requiring training or experience with that technology. Column 2 reports the respective share of technology job postings in percent which require training. To calculate shares in column 2, we only consider technology bigrams that appear in at least 100 earnings calls.

Appendix Table 15 – Top industries by technology

Emergence Year (1)	Technology (2)	Top Pioneer Industry (3)
1979	Hard disk drive	Computer and Peripheral Equipment Manufacturing
1980	Barcode reader	Computer and Peripheral Equipment Manufacturing
1981	Laser diode	Computer and Peripheral Equipment Manufacturing
1982	Personal computer	Computer and Peripheral Equipment Manufacturing
1983	Flatpanel display	Semiconductor and Other Electronic Component Manufacturing
1984	User interface	Computer and Peripheral Equipment Manufacturing
1985	Mobile phone	Communications Equipment Manufacturing
1986	Facial recognition system	Software Publishers
1987	Digital video	Computer and Peripheral Equipment Manufacturing
1988	Model organism	Pharmaceutical and Medicine Manufacturing
1989	Mobile device	Computer and Peripheral Equipment Manufacturing
1990	Debit card	Management, Scientific, and Technical Consulting Services
1991	Flash memory	Computer and Peripheral Equipment Manufacturing
1992	Machine learning	Software Publishers
1993	Financial instrument	Computer and Peripheral Equipment Manufacturing
1994	Active users	Management, Scientific, and Technical Consulting Services
1995	Hybrid electric vehicle	Motor Vehicle Manufacturing
1996	Digital content	Software Publishers
1997	Multicore processor	Semiconductor and Other Electronic Component Manufacturing
1998	Information privacy	Software Publishers
1999	Unmanned aerial vehicle	Aerospace Product and Parts Manufacturing
2000	Transaction account	Non-depository Credit Intermediation
2001	Smartphone	Information Services
2002	Online game	Software Publishers
2003	Social networking service	Other Information Services
2004	Electronic discovery	Pharmaceutical and Medicine Manufacturing
2005	LED circuit	Semiconductor and Other Electronic Component Manufacturing
2006	Augmented reality	Semiconductor and Other Electronic Component Manufacturing
2007	Self-driving car	Other Information Services

**Notes:** This table lists the top pioneer industry (NAICS four-digit codes) (in column 3) by percentage of job postings associated with a sample of technologies (in column 2) corresponding to emergence years (in column 1). The list of technologies is the same as in Table 1.

Appendix Table 16 – Wikipedia titles and trigrams

Wikipedia Title (1)	Trigrams (2)	Wikipedia Title (1)	Trigrams (2)
Web-RTC	real time communications	DNA-binding domain	dna binding domain
Real-time operating system	time operating system; real time operating	LTE telecommunication	term evolution lte
Injection molding machine	injection molding machine	Perchloric acid	sulfuric acid nitric
MEMS	electro mechanical systems; electro mechanical system	Phase detector	phase frequency detector
Single-mode optical fiber	single mode fiber; single mode fibers	Vinyl-sulfonic acid	meth acrylic acid
Heat recovery steam generator	waste gas stream; heat recovery steam	Water-gas shift reaction	water gas shift
Mobile phones on aircraft	portable electronic devices	Calcium oxide	calcium oxide calcium
Optical fiber connector	optical fiber connector; fiber optic connectors; fiber optic connector	Carboxylic acid	alkyl alkenyl aryl
Colony-stimulating factor	colony stimulating factors	Cartesian coordinate system	cartesian coordinate system
Granulocyte colony-stimulating factor	colony stimulating factor	Cellulose acetate phthalate	cellulose acetate phthalate
Protein sequencing	amino acid sequences; terminal amino acid; terminal amino acids	Complement-dependent cytotoxicity	dependent cytotoxicity cdc
Carbon-fiber-reinforced polymers	carbon fiber reinforced	Dent corn	hybrid corn variety
Transaction processing system	transaction processing system	Dibenzylideneacetone	pale yellow solid
History of mobile phones	mobile communication devices	Dimethyl sulfoxide	dimethyl sulfoxide dmso
Engine control unit	engine control unit	Fluorinated ethylene propylene	fluorinated ethylene propylene
Hollow fiber membrane	hollow fiber membrane; hollow fiber membranes	Glyceraldehyde 3phosphate dehydrogenase	glyceraldehyde phosphate dehydrogenase
Satellite navigation	navigation satellite system	Halide	fluoride chloride bromide
Polyvinyl alcohol	poly vinyl alcohol	Ion chromatography	cation exchange chromatography
Electronic control unit	electronic control module	Machine-readable medium and data	machine readable media; computer readable medium; computer readable media; machine readable medium
USB mass storage device class	mass storage devices	Magnesium sulfate	white crystalline solid
Tunnel magnetoresistance	magnetic tunnel junction	Nitrobenzene	pale yellow oil
Ultra-high-molecular-weight polyethylene	molecular weight polyethylene	Phase-shift keying	phase shift keying
Polymethyl methacrylate	poly methyl methacrylate	Phase-transfer catalyst	phase transfer catalyst
Amino acid	amino acid residues	Plant tissue culture	plant cell tissue
Attribute-based access control	access control policy	Sperm-mediated gene transfer	mediated gene transfer
Protein kinase inhibitor	protein kinase inhibitors	Zinc selenide	light yellow solid
Carboxymethyl cellulose	sodium carboxymethyl cellulose		

**Notes:** This table reports technology trigrams (in column 2) and associated technologies/Wikipedia titles (in column 1). From this list, we exclude trigrams that contain one of the baseline technology bigrams or that correspond to a Wikipedia title (technology) already associated with a technology bigram. Entries are sorted by mentions in earnings calls.

Appendix Table 17 - Wikipedia titles and unigrams

Wikipedia title	Unigrams	Wikipedia title	Unigrams	Wikipedia title	Unigrams
(1)	(2)	(1)	(2)	(1)	(2)
Website	websites	Tomography	tomography	Biomass	biomass
HTTPS	https	Flash ADC	adcs	Multiplexer	multiplexer
Internets	internets	RTP payload formats	rtp	Nonlinear system	nonlinear
Infrastructure	infrastructures	USB	usb	Multiplexing	multiplexes
Video	videos	Biomarker	biomarkers	Playlist	playlists
Upload	uploading ; upload	Defibrillation	defibrillators	Microarray	microarray
PDF	pdf	Pixel	pixels	Transgene	transgenic
Optimizing compiler	optimizations	Apheresis	apheresis	Gac	gac
Blog	blogs	SUV	suv	Zigbee	zigbee
Middleware	middleware	Bluetooth	bluetooth	Plasmid	plasmid
Intranet	intranets	Solidstate drive	ssds	Waypoint	waypoints
PMOS logic	pmos	Gallium arsenide	gaas	Clear aligners	aligner
Nondeliverable forward	dfs	Pulse oximetry	oximetry	Messenger RNA	mrnas
Prototype	prototyping	Backpack	backpacks	Cannula	cannulas
Endoscopy	endoscopy	Landline	landline	Radiation therapy	radiotherapy
United States Postal Service	usps	PowerPC	rpc	Camcorder	camcorder
Encryption	encrypt ; encryption	Hyperlink	hyperlink	Complex programmable logic device	plds
List of datasets for machinelearning research	datasets	Photonics	photonics	Liquid crystal on silicon	lcos
In vitro fertilization	vfs	Interventional radiology	endovascular	Excipient	excipients
CpG site	cpg	Discrete Fourier transform	dft	OnStar	onstar
Server Message Block	cifs ; smb	Gameplay	gameplay	Histogram	histogram
Polychlorinated biphenyl	pcbs	PCI Express	pcie	Polymerase	polymerases
Pharmacogenomics	genomics	Proteomics	proteomic	Biosensor	biosensor
Billable hours	billable	Integrated circuit	microchips	Hydrocodone	hydrocodone
Topology	topologies	Television timeout	timeout	Oxymorphone	oxymorphone
Clinical decision support system	dsss	Floating production storage and offloading	offloading	TFT LCD	tfts
Interoperability	interoperability	Satellite navigation	gnss	Insulated gate bipolar transistor	igbts
Trusted Platform Module	tpm	Nebulizer	nebulizers	Xenotransplantation	xenograft
Biofuel	biofuel	Biomaterial	biomaterials	Antibody	immunoglobulin
Realtime operating system	rtos	Multimodality	multimodal	Sunscreen	sunscreen

Wikipedia title (1)	Unigrams (2)	Wikipedia title (1)	Unigrams (2)
Linear particle accelerator	linac	Erlotinib	erlotinib
Rituximab	rituximab ; rituxan	Steam assisted gravity drainage	sagd
Breakpoint	breakpoint	Moisturizer	moisturizing
Embolization	embolization	Supercapacitor	ultracapacitor
Collagen	collagens	Protein isoform	isoforms
Reticle	reticle	Stevia	stevia
Through silicon via	tsvs	Paclitaxel	paclitaxel
Penicillin	penicillins	PSMA scan	psma
Dopamine	dopamine	Lysine	lysines
Duplexer	duplexer	Doxorubicin	doxorubicin
Syringe	prefilled	Prostaglandin	prostaglandins
Megabyte	megabytes	Calcitonin	calcitonin
Small Formfactor		Tamoxifen	tamoxifen
Pluggable	pluggable	Irinotecan	irinotecan
Glycemic index	glycemic	Bupropion	bupropion
Naloxone	naloxone	Cyclophosphamide	cyclophosphamide
Photomask	photomask	Atorvastatin	atorvastatin
Angiogenesis	angiogenic; angiogenesis	Budesonide	budesonide
Microparticle	microsphere	Cytarabine	cytarabine
Biopolymer	biopolymer	Sulfonylurea	sulfonylureas
Revascularization	revascularization		
Joystick	joysticks		
Glycoprotein	glycoprotein		
Backlight	backlights		
Carbonless copy paper	carbonless		
Cephalosporin	cephalosporins		
Polymethyl methacrylate	pmma		
Neurostimulation	neurostimulation		
Chemokine	chemokines		
Cetuximab	erbitux		
Nucleoside	nucleosides		

**Notes:** This table reports the unigrams (in column 2) and their respective technologies/Wikipedia titles (in column 1), in the order of total job postings that they appear in. The table lists unigrams that appear in at least 100 earnings calls. Out of a total of 200 such unigrams, we exclude 53 unigrams that correspond to existing Wikipedia titles (technologies) already associated with a technology bigram.

Appendix Table 18 – Comparison of main results: Disruptive and non-disruptive bigrams

Panel A: Disruptive bigrams				
Dependent Variable	Share College Educated	Coefficient of Variation	Normalized Share	Log (Coefficient of Variation)
Result:	Skill Broadening	Region Broadening	Pioneer Persistence	Differential Skill Broadening
	(1)	(2)	(3)	(4)
Year since emergence	<b>-0.541***</b> (0.046)	<b>-0.147***</b> (0.015)		-0.033*** (0.003)
Pioneer			1.084*** (0.309)	
Pioneer * Year since emg.			-0.032** (0.013)	
Low skill * Year since emg.				<b>-0.009***</b> (0.003)
Constant	64.289*** (1.069)	7.400*** (0.339)	0.621*** (0.001)	2.726*** (0.061)
R-squared	0.895	0.816	0.038	0.852
N	4,270	4,270	3,965,122	12,706
Fixed Effects	Bigram	Bigram	Bigram, CBSA, Year	Bigram
Technologies	428	428	428	428
Rate of decline per year	0.008	0.020	<b>-0.029 (0.005)</b>	
Years to Full Decay	62.87	50.43	33.88	.
Panel B: Intermediate and non-disruptive bigrams				
Result:	Skill Broadening	Region Broadening	Pioneer Persistence	Region Broadening By Skill
	(1)	(2)	(3)	(4)
Year since emergence	<b>-0.413***</b> (0.053)	<b>-0.164***</b> (0.017)		-0.022*** (0.002)
Pioneer			1.986*** (0.392)	
Pioneer * Year since emg.			-0.051*** (0.015)	
Low skill * Year since emg.				<b>-0.007**</b> (0.003)
Constant	63.316*** (1.273)	9.977*** (0.408)	0.582*** (0.001)	2.910*** (0.046)
R-squared	0.919	0.807	0.023	0.796
N	4,307	4,307	3,999,480	12,784
Fixed Effects	Bigram	Bigram	Bigram, CBSA, Year	Bigram
Technologies	431	431	431	431
Rate of decline per year	0.007	-0.016	<b>-0.026 (0.003)</b>	.
Years to Full Decay	79.95	60.84	38.94	.
Difference vs disruptive	-0.128(0.070) *	0.017(0.022)	-0.003(0.006)	-0.002(0.005)

**Notes:** This table reports results from regressions corresponding to our primary results: skill broadening (column 1), region broadening (column 2), pioneer persistence (column 3), and differential region broadening by skill (column 4), separately for disruptive (Panel A) and other (intermediate + non-disruptive, Panel B) technologies. Our primary coefficients for comparison are in bold in both panels. Panel A is restricted to the sample of technology bigrams that appear in at least 100 earnings calls. Panel B uses all other technology bigrams that appear in at least 1000 job postings in our sample. For more details on the specification, refer to the notes in Table 4, 5, 6, and 7.

Appendix Table 19: Spread of high vs. low-skill jobs, using all bigrams

Panel A: Region Broadening			
	$\log(\text{Coefficient of Variation})_{\tau,t}$		
	(1)	(2)	(3)
Skill Level	High	Medium	Low
$Years\ since\ emergence_{\tau,t}$	-0.024*** (0.002)	-0.030*** (0.002)	-0.035*** (0.002)
R-squared	0.842	0.853	0.855
N	8,335	8,330	8,219
Bigrams	835	835	835
Bigram FE	YES	YES	YES
Std. Errors (cluster)	Wiki Title	Wiki Title	Wiki Title
Panel B: Pioneer Advantage			
	$Normalized\ Share_{c,\tau,t}$		
	(1)	(2)	(3)
Skill Level	High	Medium	Low
$Pioneer_{c,\tau}$	1.309*** (0.233)	0.909*** (0.225)	1.104*** (0.253)
$Pioneer_{c,\tau}$ * $year\ since\ emg_{\tau,t}$	-0.029*** (0.010)	-0.022** (0.010)	-0.036*** (0.010)
R-squared	0.016	0.018	0.012
N	8,491,998	8,478,284	8,319,989
Bigrams	835	835	835
Bigram FE	YES	YES	YES
Std. Errors (cluster)	Wiki Title	Wiki Title	Wiki Title
<b><i>Rate of decline per year</i></b>	<b>-0.022</b> <b>(0.004)</b>	<b>-0.024</b> <b>(0.005)</b>	<b>-0.032</b> <b>(0.003)</b>

**Notes:** This table reports region-broadening regressions at the technology bigram x year level (Panel A) and pioneer advantage regressions at the technology bigram x year x location level (Panel B), now using all technology bigrams with at least 1000 job postings. These analyses are akin to those in columns 1 and 2 of Table 6, Panels A and B, but now include all technology bigrams that appear in at least 1000 job postings. Columns 1, 2, and 3 show separate regressions for high-skill, medium-skill, and low-skill job postings. Panel A reports results from regressions of  $\log(\text{Coefficient of Variation})_{\tau,t}$  on the number of years since the technology's year of emergence. Panel B reports results from regressions of the  $Normalized\ share_{c,\tau,t}$  on a dummy indicating pioneer status of the CBSA and the interaction of this dummy with number of years that have elapsed since bigram's emergence. The rate of decline per year is calculated as  $\frac{\beta_D}{\beta_P}$ , where  $\beta_P$  is the coefficient of  $Pioneer_{c,\tau}$  and  $\beta_D$  is the coefficient of  $Pioneer_{c,\tau} * Years\ since\ emg_{\tau,t}$ . For more details on the specification, refer to the note in Table 6.



# Appendix Figure 1 – Examples of patent text



US007406200B1

(12) **United States Patent**  
**Syeda-Mahmood et al.**

(10) **Patent No.:** **US 7,406,200 B1**

(45) **Date of Patent:** **Jul. 29, 2008**

(54) **METHOD AND SYSTEM FOR FINDING STRUCTURES IN MULTI-DIMENSIONAL SPACES USING IMAGE-GUIDED CLUSTERING**  
 2004/0202370 A1 10/2004 Fisher et al.  
 2007/0185946 A1\* 8/2007 Basri et al. .... 708/200

(75) **Inventors:** **Tanveer Syeda-Mahmood**, Cupertino, CA (US); **Peter J. Haas**, San Jose, CA (US); **John M. Lake**, Cary, NC (US); **Guy M. Lohman**, San Jose, CA (US)

**OTHER PUBLICATIONS**

Gdalyhu et al., "Self-Organization in Vision: Stochastic Clustering for Image Segmentation, Perceptual Grouping, and Image Database Organization", IEEE Trans. on Pattern Analysis and Machine Intelligence, Oct. 2001, pp. 1053-1074, vol. 23, No. 10.

(73) **Assignee:** **International Business Machines Corporation**, Armonk, NY (US)

\* cited by examiner

(\* ) **Notice:** Subject to any disclaimer, the term of this patent is extended or adjusted under 35 U.S.C. 154(b) by 0 days.

*Primary Examiner*—Aaron W Carter  
 (74) *Attorney, Agent, or Firm*—Cantor Colburn LLP; Van Ngyen

(21) **Appl. No.:** **11/970,946**

(57) **ABSTRACT**

(22) **Filed:** **Jan. 8, 2008**

(51) **Int. Cl.**  
**G06K 9/62** (2006.01)  
**G06F 7/00** (2006.01)  
**G06F 17/30** (2006.01)  
**G06F 17/00** (2006.01)

A method is provided clustering data points in a multidimensional dataset in a multidimensional image space that comprises generating a multidimensional image from the multidimensional dataset; generating a pyramid of multidimensional images having varying resolution levels by successively performing a pyramidal sub-sampling of the multidimensional image; identifying data clusters at each resolution level of the pyramid by applying a set of perceptual grouping constraints; and determining levels of a clustering hierarchy by identifying each salient bend in a variation curve of a magnitude of identified data clusters as a function of pyramid resolution level.

(52) **U.S. Cl.** ..... **382/225**; 707/3; 707/104.1

(58) **Field of Classification Search** ..... 382/224–228; 707/3–6, 104.1

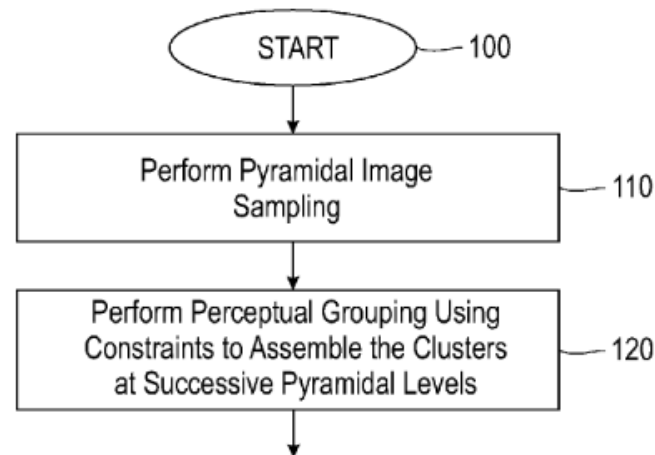
See application file for complete search history.

(56) **References Cited**

U.S. PATENT DOCUMENTS

2004/0013305 A1 1/2004 Brandt et al.

**5 Claims, 10 Drawing Sheets**



US 7,406,200 B1

**I**  
**METHOD AND SYSTEM FOR FINDING STRUCTURES IN MULTI-DIMENSIONAL SPACES USING IMAGE-GUIDED CLUSTERING**

**CROSS REFERENCE TO RELATED APPLICATIONS**

The present application is co-pending with the concurrently filed application, entitled "METHOD AND SYSTEM FOR DATA CLASSIFICATION BY KERNEL DENSITY SHAPE INTERPOLATION OF CLUSTERS," assigned to the assignee of the present application, the contents of which are incorporated herein by reference in their entirety.

**BACKGROUND OF THE INVENTION**

1. **Field of the Invention**  
 This invention relates to data clustering, and more particularly, to the clustering of multidimensional data to determine high-level structures.

2. **Description of Background**  
 Data clustering (or just clustering) is the categorization of objects into different groups, or more precisely, the organizing of a collection of data into clusters, or subsets, based on quantitative information provided by one or more traits or characteristics shared by the data in each cluster. A cluster is a collection of objects which are "similar" between them and "dissimilar" to the objects belonging to other clusters. The goal of clustering is to determine an intrinsic grouping, or structure, in a set of unlabeled data. For example, the functional dependency between two or more time series can lie along a curve. As an example, FIG. 1 shows a graph of a functional dependency between a pair of time series that maps to a perceptible curve having a rotated U-like structure. Clustering can be used to perform statistical data analysis in many fields, including machine learning, data mining, pattern recognition, medical imaging and other image analysis, and bioinformatics.

For applications dealing with sets of high-dimensional data such as multimedia processing applications (for example, content-based image and video retrieval, multimedia browsing, and multimedia transmission over networks), the finding of underlying high-level structures by clustering and categorization is a fundamental analysis operation. A good clustering scheme should, for example, help to provide an efficient organization of content, as well as provide for better retrieval based upon semantic qualities. In video retrieval, because of the larger number of additional features resulting from motion in time, efficient organization is particularly important. In image-based retrieval, semantic quality retrieval is particularly important because clustering provides a means for grouping images into classes that share some common semantics.

Even though clustering of multidimensional datasets is important to determining high-level structures, much of the focus in multidimensional data analysis has been on feature extraction and representation, and existing methods available from data mining and machine learning have been relied on for the clustering task. These methods are primarily based upon the similarity criterion of distance or proximity in which two or more objects belong to the same cluster if they are "close" according to a given distance function that defines a distance between elements of a set (for example, the simple Euclidean distance metric).

The nature of multidimensional datasets, however, presents a number of peculiarities that can lead to misleading or

insufficient results using distance-based clustering, particularly for cases of grouping high-dimensional objects into high-level structures. First, the number of feature dimensions in multidimensional datasets tends to be large in comparison to the number of data samples. As an example, a single four second action video assuming a pair of features per frame (for instance, for representing the motion of the object centroid) can have at least 240 feature dimensions. Similarly, in image clustering, while color, texture, and shape features can encompass hundreds of features, the number of samples available for training could be comparably small. This can result in a data space that is high-dimensional but sparse. The sparseness of the data points can make it difficult to identify the clusters because observation at multiple scales may be needed to spot the patterns.

A second issue that may arise is that the number of clusters for a multidimensional dataset is often unknown and more than one set of clusters may be possible. Different relative scalings can lead to groupings with different structures, even with measurements being taken in the same physical units. To make an informed decision as to relative scaling using existing clustering methods, either the number of clusters needs to be known a priori or a hierarchical clustering must be performed that yields several possible clusters without a specific recommendation on one. In a hierarchical clustering, the process builds (agglomerative), or breaks up (divisive), a hierarchy of clusters. The traditional representation of such a hierarchy of clusters is a tree structure called a dendrogram, which depicts the mergers or divisions which have been made at successive levels in the clustering process. A bottom row of leaf nodes represent data and the set of remaining nodes represent the clusters to which the data belong at each successive stage of analysis. The leaf nodes are spaced evenly along the horizontal axis, and the vertical axis gives the distance (or dissimilarity measure) at which any two clusters are joined. Divisive methods begin at the top of the tree, while agglomerative methods begin at the bottom, and cutting the tree at a given height will give a clustering at a selected precision. The bottom level of the hierarchy includes all data points as one cluster.

As an example of the scaling issue, a clustering scenario is provided that involves a type of dataset for which the structure of the functional dependency between two or more time series can take a variety of forms. As an example, FIG. 2 illustrates a graph of functional dependencies between a pair of time series in which the noticeable structures are that of three separate lines radiating from common points. While different structures from within this graph may be obtained using hierarchical clustering methods, ideally, it would be desirable to have the result of clustering the dataset indicate the lower level structures (such as the individual splotches in FIG. 2) as well as the higher-level structures formed (such as the lines perceived in FIG. 2) without necessarily leading to a single cluster at the top level, unless that is in fact matching how the data collection should be perceived.

**SUMMARY OF THE INVENTION**

The shortcomings of the prior art can be overcome and additional advantages can be provided through exemplary embodiments of the present invention that are related to method for clustering data points in a multidimensional dataset in a multidimensional image space. The method comprises generating a multidimensional image from the multidimensional dataset; generating a pyramid of multidimensional images having varying resolution levels by successively performing a pyramidal sub-sampling of the

**Notes:** These figures show examples of patent text. The left panel displays the first page of U.S. Patent 7406200, while the right panel shows the second page of text in that same patent. For analysis in our paper, we use all text sections, including title, abstract, background, summary, description, and claims.

Appendix Figure 2 – Examples of Wikipedia entries and procedure for filtering technology bigrams

The figure displays two screenshots of Wikipedia pages. The top screenshot is for the article "Mobile device". It features a navigation sidebar on the left with options like "Main page", "Contents", and "Tools". A table of contents is visible, with the first six items circled in red: "1 Characteristics", "2 Types", "3 Uses", "4 See also", "5 References", and "6 Sources". A notice at the top of the article body states: "This article needs additional citations for verification. Please help improve this article by adding citations to reliable sources. Unsourced material may be challenged and removed." Below this, the article text begins: "A mobile device (or handheld computer) is a computer small enough to hold and operate in the hand. Typically, any handheld computer device will have an LCD or OLED flat screen interface, providing a touchscreen interface..."

The bottom screenshot is for the article "Machine learning". It has a similar layout. Its table of contents is also circled in red, with the seventh item, "7 Model assessments", specifically highlighted. The article text starts with: "Machine learning (ML) is a field of inquiry devoted to understanding and building methods that 'learn', that is, methods that leverage data to improve performance on some set of tasks. It is seen as a part of artificial intelligence. Machine learning algorithms build a model based on sample data, known as training data, in order to make predictions or decisions without being explicitly programmed to do so." On the right side of the page, there are two related article boxes: "Part of a series on Machine learning and data mining" and "Part of a series on Artificial intelligence".

Notes: The figure shows a sample of Wikipedia pages corresponding to two example technologies. The red circles indicate sections that are associated with selection of technology bigrams.



## Appendix Figure 3 – Example of a Burning Glass technology job posting, and corresponding location and occupation fields

### 1.1 CanonCity

Seattle

### 1.2 CanonCountry

USA

### 1.3 CanonState

WA

### 1.4 CleanJobTitle

Server Developer

### 1.5 JobDate

1/11/2010

### 1.6 JobDomain

seattle.craigslist.org

### 1.7 JobText

Date: 2010-01-08, 10:44AM PST

Reply to:

Do you like shaking up the status quo? How about shaking up the advertising industry? Advertising is fundamentally broken in that to manage campaigns across TV, online display, search, and print requires armies of people and has no accountability on the results across the mediums. Come join Lucid Commerce, which is a start up in Pioneer Square created by advertising industry veterans from Aquantive, Microsoft, Overture, Specific Media, and SpotRunner. At Lucid, we are building the first cross-medium robotic advertising agency using deep advertiser integration and **MACHINE LEARNING**. We are already running media in television and very soon online. We are looking for a self-starter that is passionate developing systems that have low latency/high scale requirements and use **MACHINE LEARNING** techniques to optimize. Work on building the display banner ad server where we need to determine the market potential of a request in less than 30ms given a set of inbound attributes about the audience. Or work on the television

system which is running daily to determine which commercial spots are predicted to yield the highest amount of revenue given a feature set of hundreds of attributes about the programming, the buying audience's demographics, other similar product's buyer attributes, etc. This person must be able to work with the data mining team in terms of contributing to model design and take it from design in R to implementation in C#/C++. This person must also be able to learn and build statistical unit testing techniques, operational model maintenance, etc.

The ideal skills are: Understanding of HTTP request (headers, cookie setting, logging, etc); OOP concepts/methodology; **MACHINE LEARNING** or statistics background; C # or Java C++; Experience with caching systems (Ehcache, Memcached, Terracotta, etc); Working knowledge of SQL; Working knowledge of multi-variate testing techniques; Advanced degree in CS, statistics, or other strong quantitative degree a plus.

In terms of campus hiring this is a fantastic opportunity to learn an industry and technology from a highly motivated team of experts and potential to join an early stage startup set to grow. There are also limited opportunities for academic publishing.

**Location:** Seattle

**Compensation:** Market Rate

Principals only. Recruiters, please don't contact this job poster. Please, no phone calls about this job! Please do not contact job poster about other services, products or commercial interests.

PostingID: 1543312019 Copyright 2010 craigslist, inc.

### 1.8 JobID

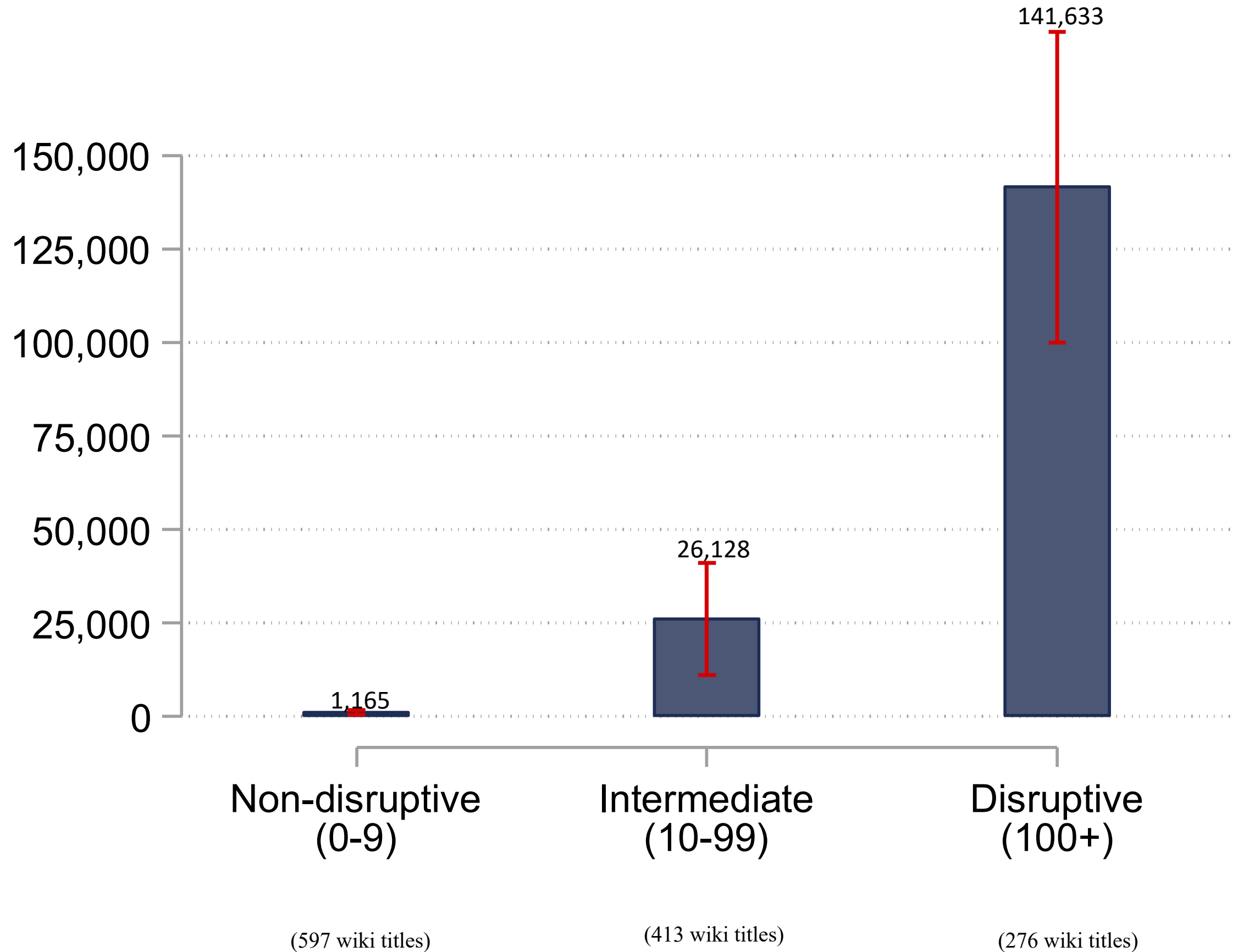
311172646

### 1.9 JobURL

<http://seattle.craigslist.org/see/sof/1543312019.html>

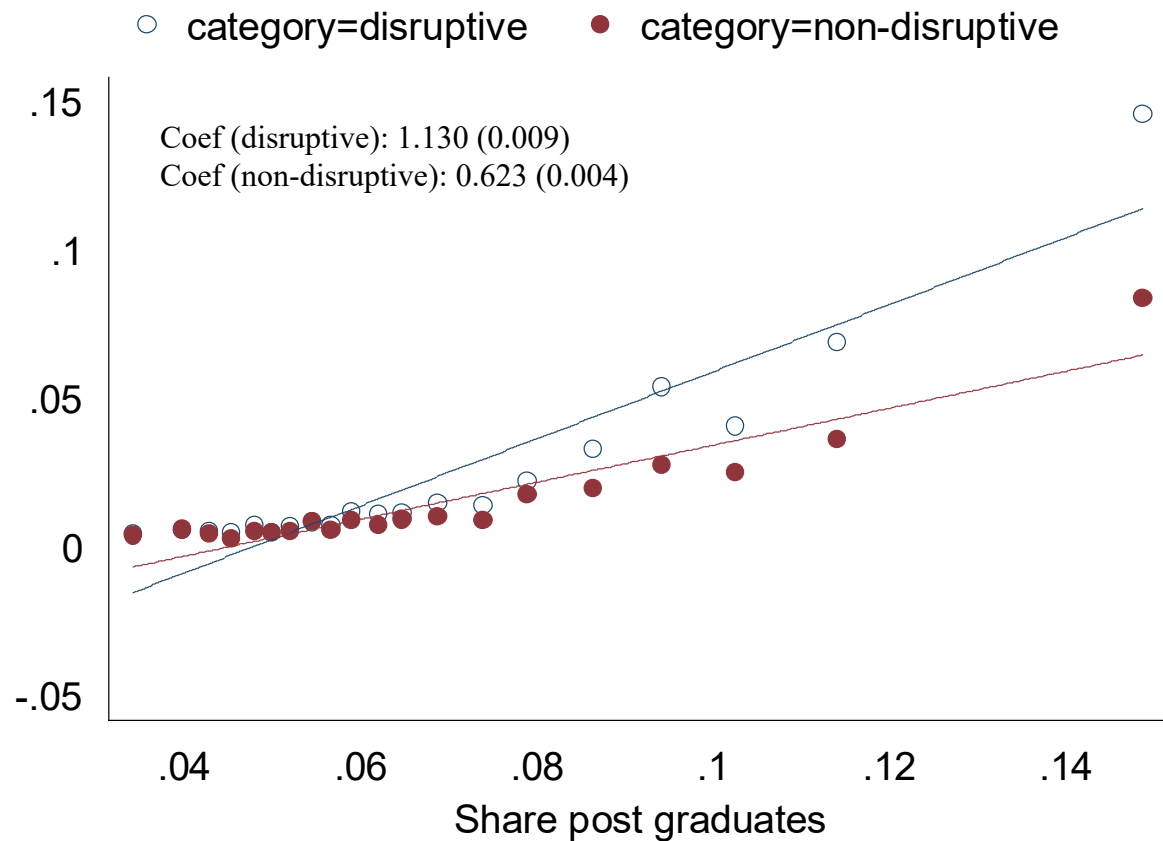
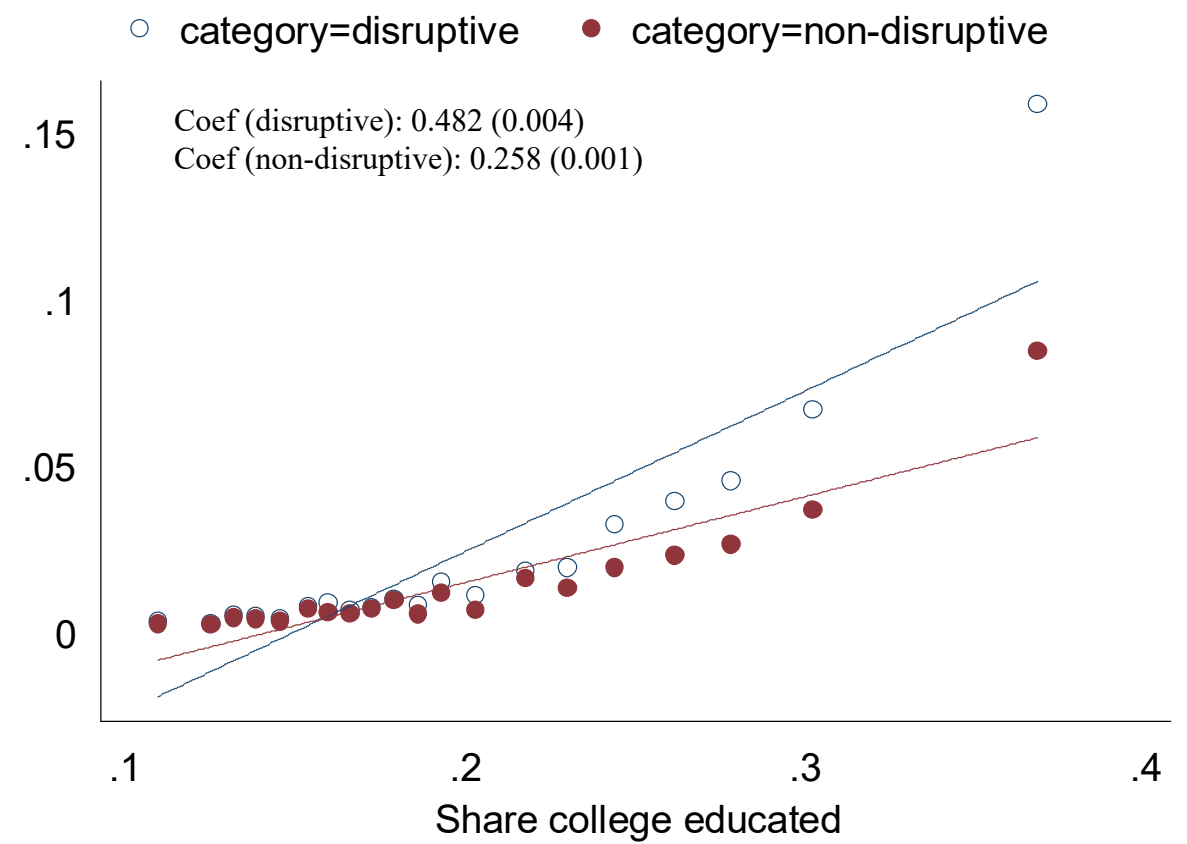
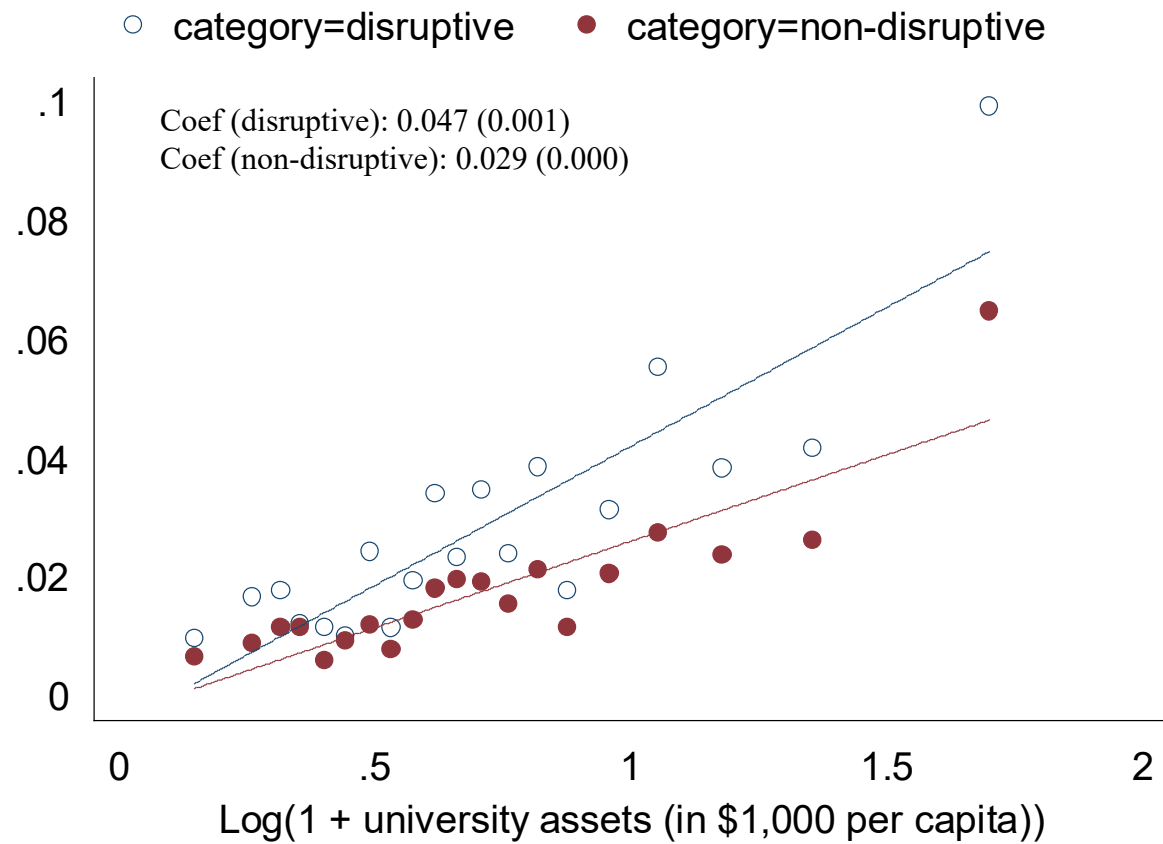
**Notes:** The figure shows an example of a job posting in Burning Glass, with associated data. This posting mentions one of our technologies (“machine learning”).

**Appendix Figure 4 – Average number of job postings by type of technology – disruptive, intermediate and non-disruptive**



**Notes:** The figure shows the average number of job postings associated with different types of technologies (disruptive, intermediate, and non-disruptive). Disruptive technologies are technologies associated with at least one bigram that appears in at least 100 earnings calls; intermediate technologies are those technologies associated with at least one bigram that appears in at least 10 earnings calls and are not one of the disruptive technologies. The rest are classified as non-disruptive technologies. Whiskers show 95%-confidence interval. A t-test of difference between the average number of disruptive and intermediate postings results in a F-stat of 57.04 (p-value of less than 0.001).

**Appendix Figure 5 – Early disruptive and non-disruptive patenting vs. local skill composition**



**Notes:** The figures show binned scatter plots of disruptive and non-disruptive patents in the ten years prior to the emergence date per 1,000 people in a CBSA over measures of skill and university presence in the CBSA. The coefficients and robust standard errors displayed are from the associated regressions with standard errors clustered by CBSA. In all panels, the coefficient for disruptive patents are (statistically) significantly higher than the coefficient for non-disruptive patents at 1% level.

## Appendix Figure 6 – Examples of use and research, develop, and production (RDP) in job postings

### Use: Sales Representative

We are in search of an outgoing driven and reliable individual who is looking for a part time or full time opportunity to become a brand rep for a regulated product in convenience store locations in your area you will be a key asset to the program...

Responsibilities and requirements:

- work on product displays pull product out of back stock and merchandise replenish displays as needed.
- use third channel technology on a smart device to collect crucial data engage with consumers and provide sales support/brand education to retail associates
- *reliable transportation a **smart phone** with internet access.*

### Use: Clinical Nurse

Position is responsible for creating and maintaining positive customer relations for hospice and, as appropriate, home health and other-in-home services programs throughout the assigned hospital/clinic.

Essential Duties:

- ...
- Documents in patient hospital medical record pertinent information related to discharge from hospital to home or facility/agency setting.
- *Appropriately documents activities in the Providence **electronic medical system**, tracks referrals received by nursing unit and accepted by each Providence agency.*

### Use: Ultrasound Technologist

The Ultrasound Technologist II under general direction, operates ultrasound equipment to obtain high quality ultrasound examinations of various body parts as ordered by the Physician. Prepares patients, processes images and assists Physicians as needed. Evening Shift and rotate call. Qualifications RDMS required. ARRT a plus. 3-5 years of experience preferred. *Successful completion of a formal (Accredited) educational /training program for Radiology and/or **Medical Imaging**.*

### Develop: Software Engineer

Job Description: *Analyze requirements, develop the **software architecture**, plan and perform detailed software design activities, specify software test requirements, estimate software size and development effort and plan for future product releases is required.*

### Develop: Scientist

[The company] has an opening in R&D for an individual experienced in radiological imaging to lead the development of computed tomography (CT) systems used in pre-clinical imaging research. We are emphasizing a multi-modality imaging approach, combining CT with optical molecular imaging. The successful candidate will guide the development of both hardware and software platforms, working closely with a multi-disciplinary team of engineers, software developers, and biologists...

*With a keen focus on clinically-relevant experimentation, [this company's] portfolio of offerings includes state-of-the-art microfluidics, lab automation & **liquid handling**, optical imaging technologies, and discovery & development outsourcing solutions....*

### Develop: Modeling & Simulation Software Engineer

Our client is seeking a motivated and passionate scientist/engineer with the technical vision required to translate advanced modeling and simulation concepts into operational reality. (...) The successful applicant will function as a key technical member of an interdisciplinary team in a dynamic development environment, providing leadership across the full R& D cycle. Position responsibilities will include concept formulation, proposal writing, algorithm design, software implementation and integration, client interaction, and research publication...

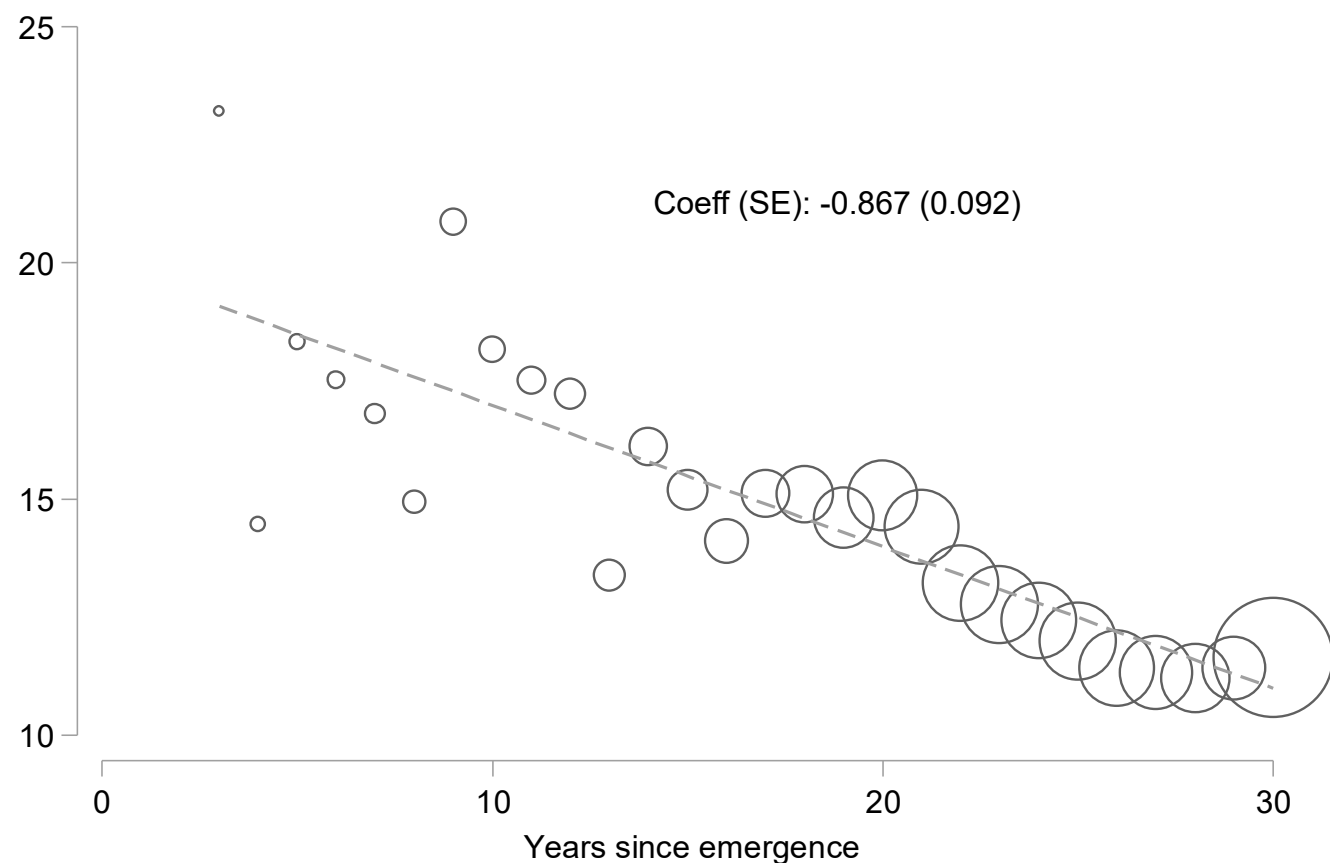
Ph.D. in Computer Science or M.S. in Computer Science plus minimum 5 years experience . Hands-on experience with one or more advanced modeling approaches, transforming theory into real world solutions. *Strong software development skills including systems integration, network programming, AI programming, and **virtual environments**.*

**Notes:** The panels above show examples of job postings mentioning our technology bigrams categorized into use, and research, development, and production.

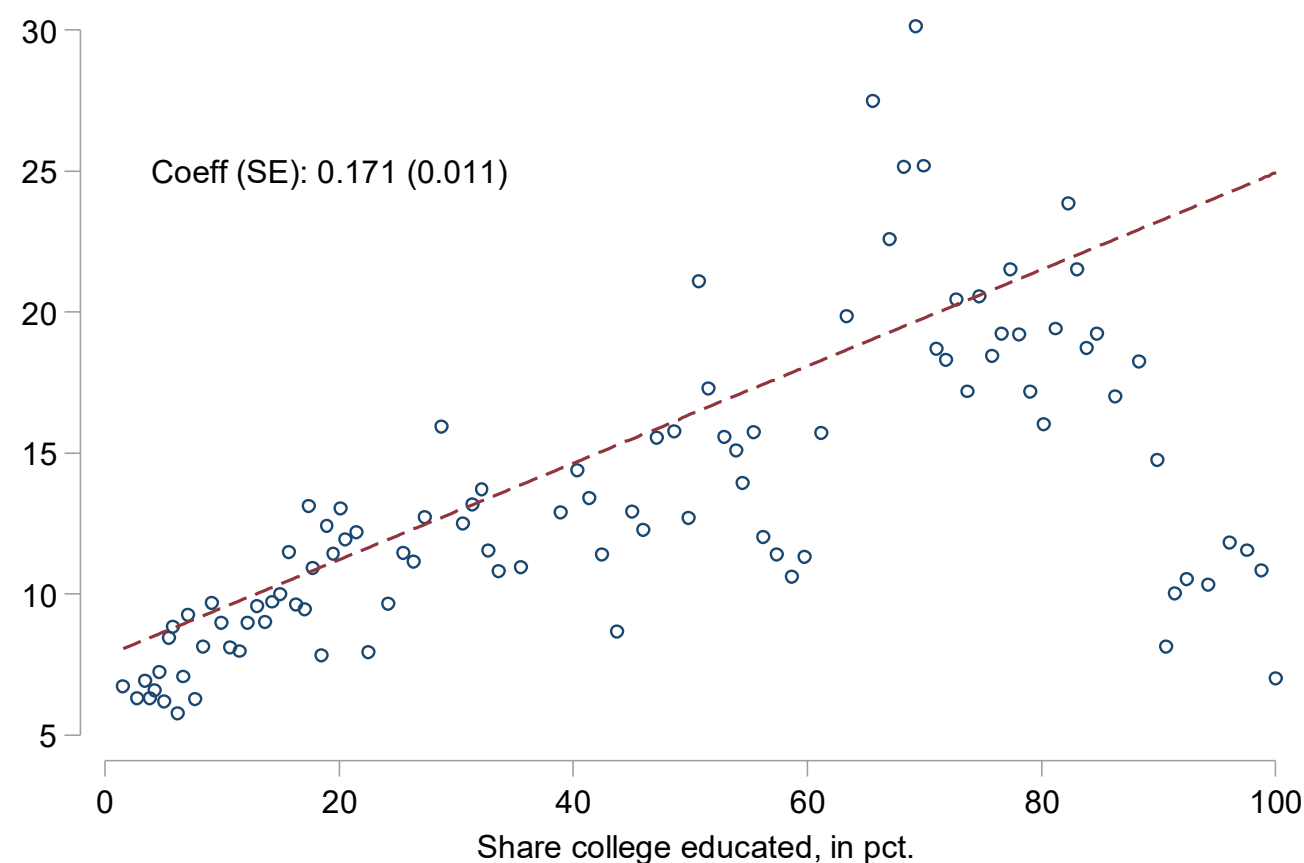


## Appendix Figure 7 – Technology job postings with RDP synonyms, by year since emergence and college education

**Panel A: Pct. of technology RDP job postings by year since emergence**



**Panel B: Pct. of technology RDP jobs by college education**

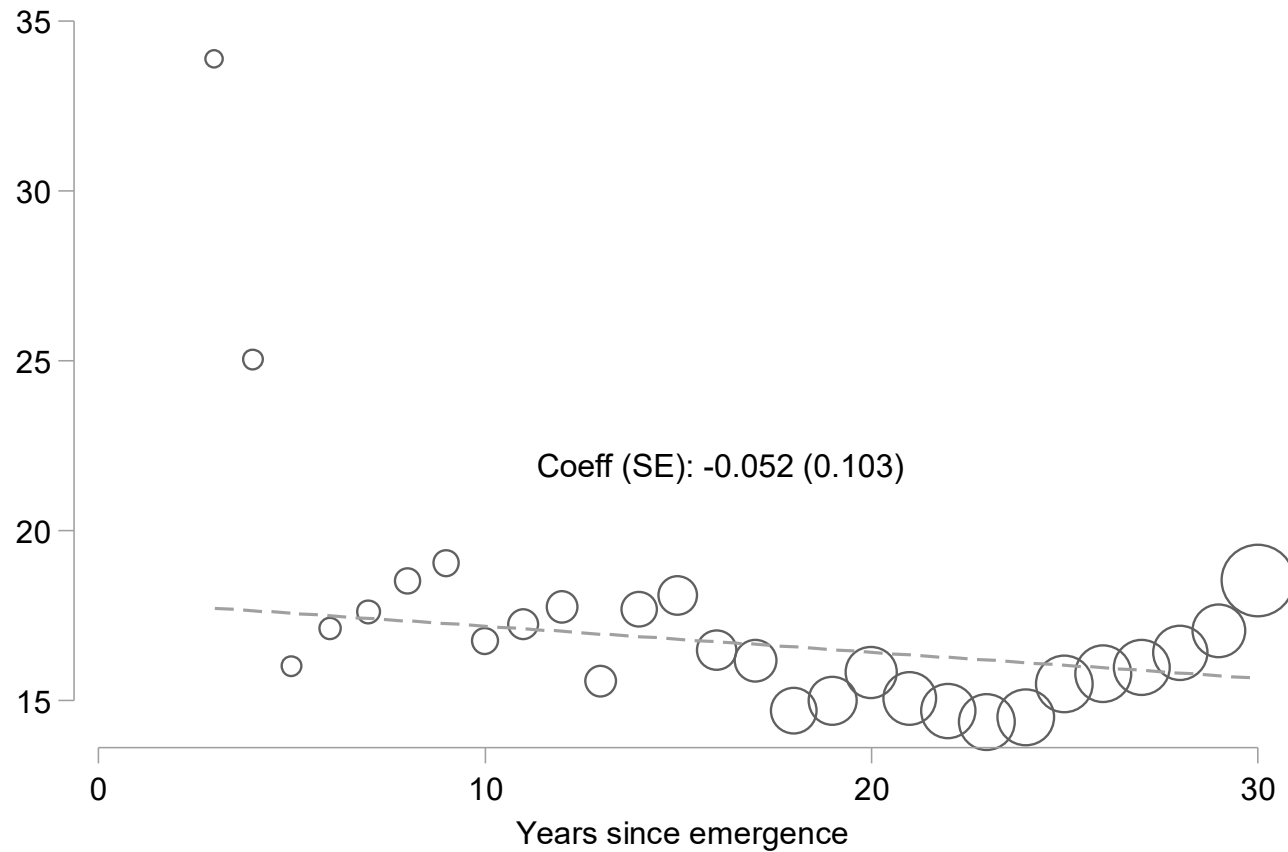


**Notes:** For a panel of technology bigram x year observations, the figure in Panel A shows a binned scatter plot of the share of Research, Development, and Production (RDP) job postings relative to the years since emergence. Each bin represents the average share of job postings categorized as RDP jobs across technology bigrams for a given year since emergence. The sample pools across technologies, where each observation in the sample denotes a technology bigram x year observation. Only observations at the year of emergence and after are included. The size of circles is proportional to the square root of the number of technology bigram job postings, capped at 100.

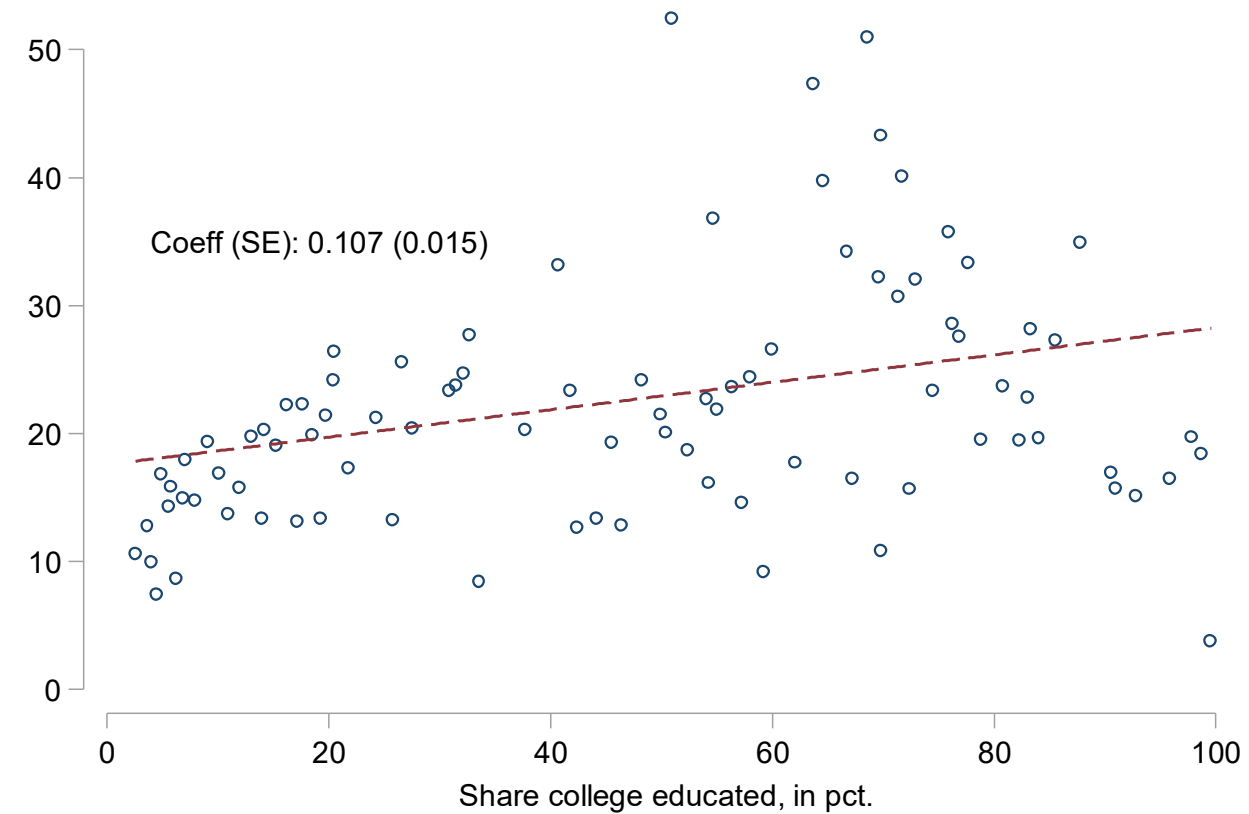
**Notes:** The figure in Panel B plots a binned scatter plot of share of Research, Development, and Production (RDP) postings against the share of college-educated postings. Each observation in the dataset is at the technology bigram x year level.

## Appendix Figure 8 – Technology job postings with training synonyms, by year since emergence and college education

**Panel A: Pct. of technology training job postings by year since emergence**



**Panel B: Pct. of technology training jobs by college education**



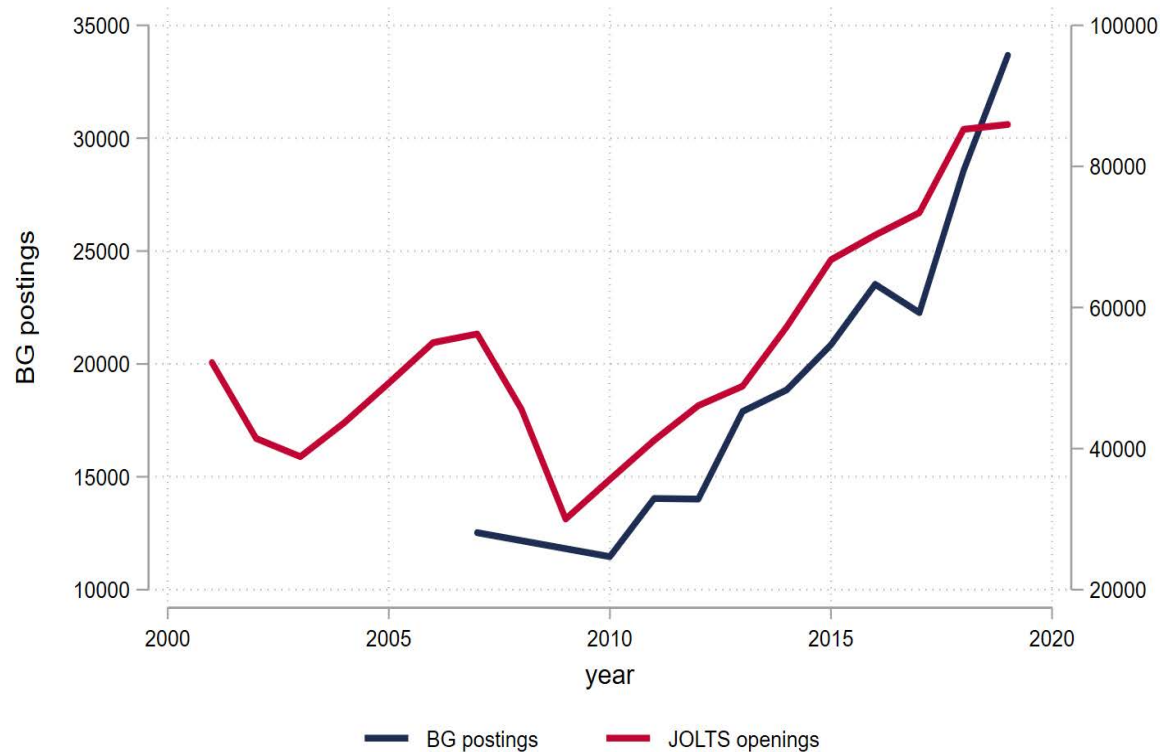
**Notes:** For a panel of technology bigram x year observations, the figure in Panel A shows a binned scatter plot of the share of training postings against the years since emergence. Each bin represents the average share of job postings categorized into 'training' across technology bigrams for a given year since emergence. The sample pools across technology bigrams, where each observation in the sample denotes a technology bigram x year observation. Only observations at the year of emergence and after are included. The size of circles is proportional to the square root of the number of technology bigram job postings, capped at 100.

**Notes:** The figure in Panel B plots a binned scatter plot of share of training postings against the share of college-educated postings. Each observation in the dataset is at the technology bigram x year level.

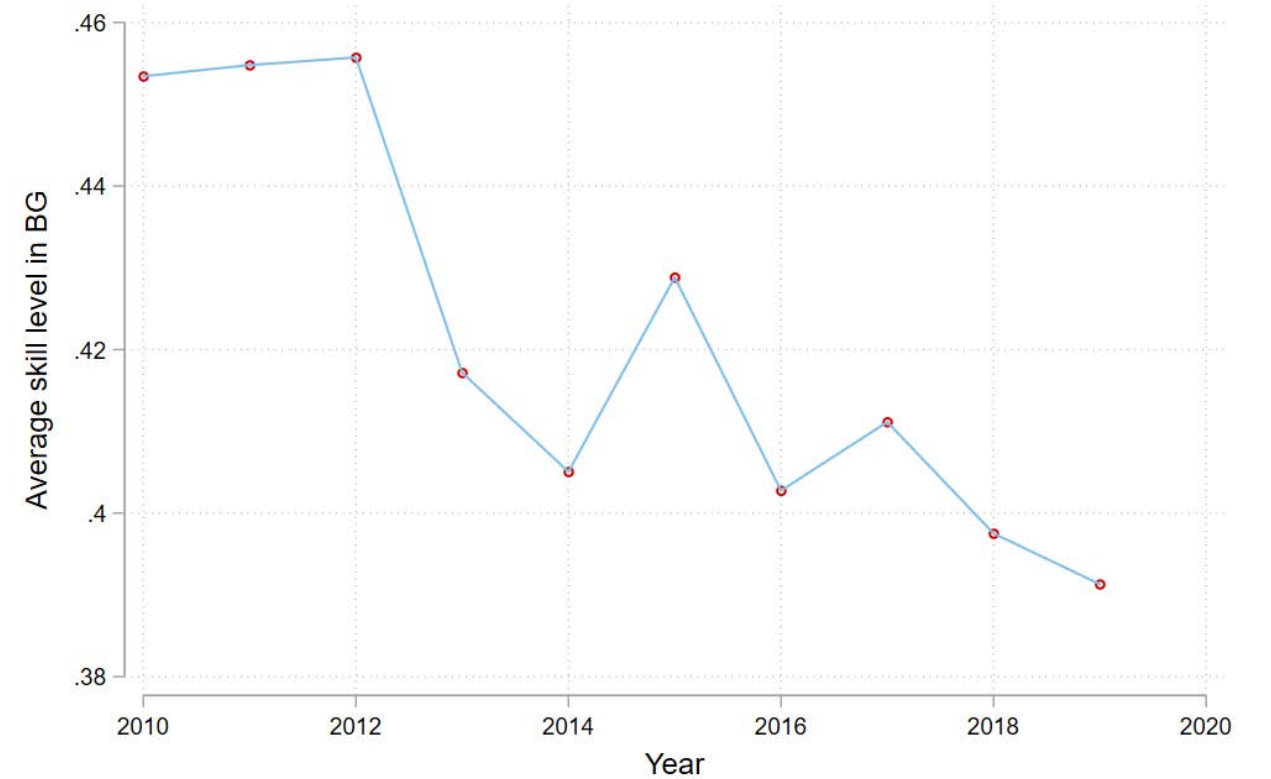


# Appendix Figure 9 – Total number and composition of Burning Glass job postings, over time

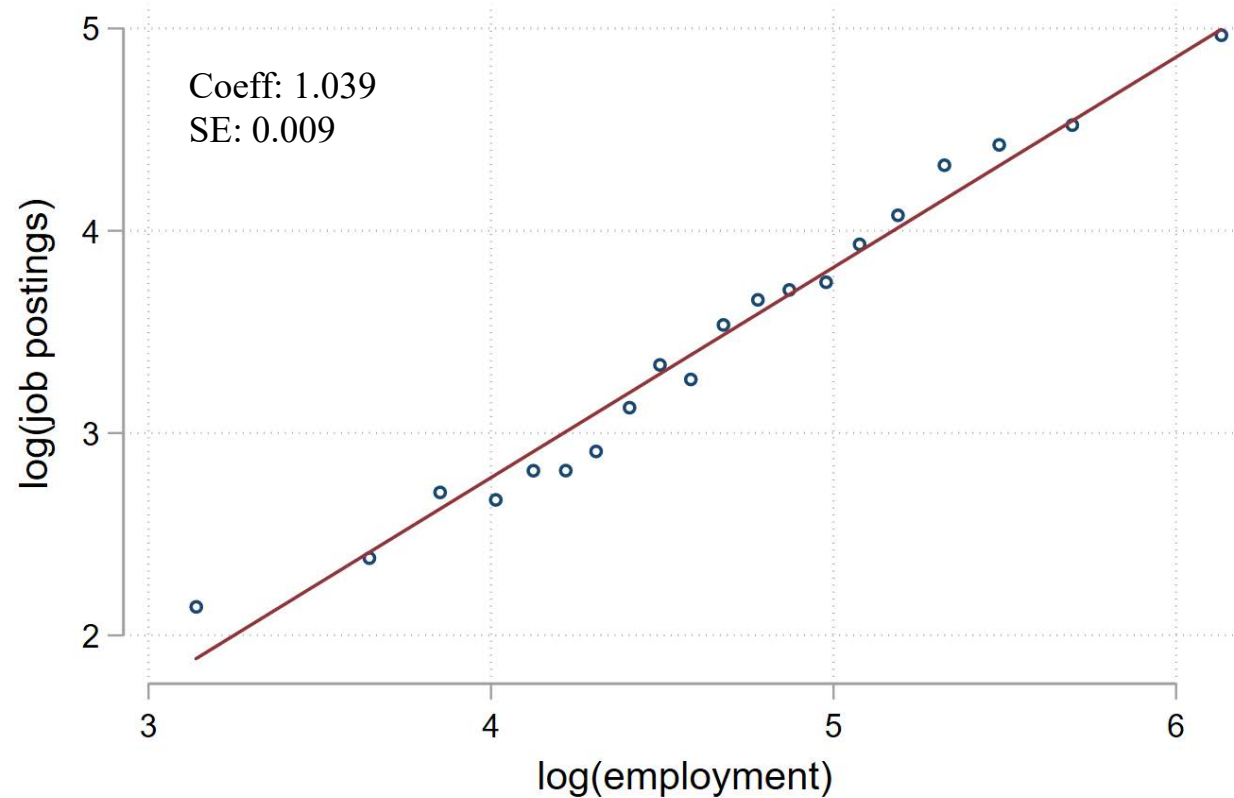
## Panel A: Number of job postings in BG vs JOLTS



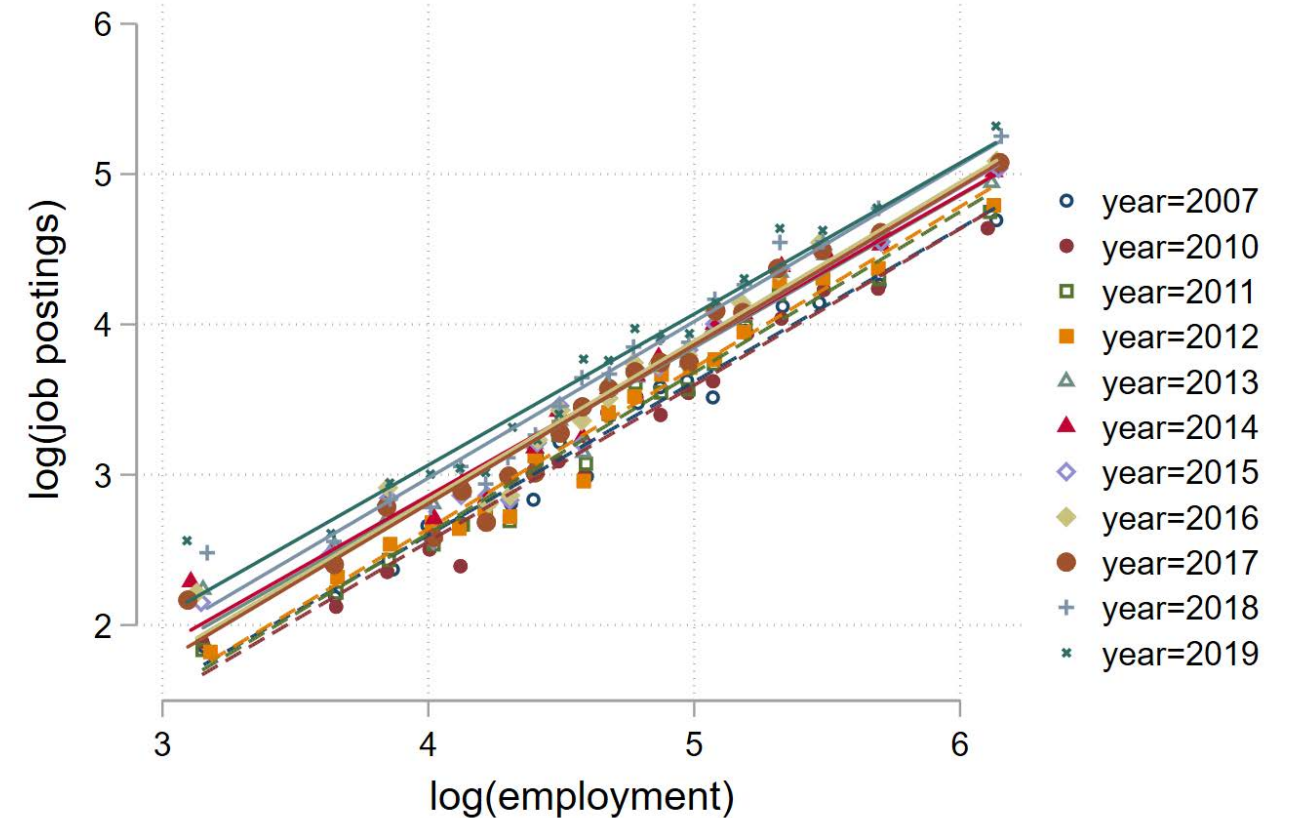
## Panel B: Share of job postings requiring a college-education



## Panel C: Number of job postings vs employment in OES

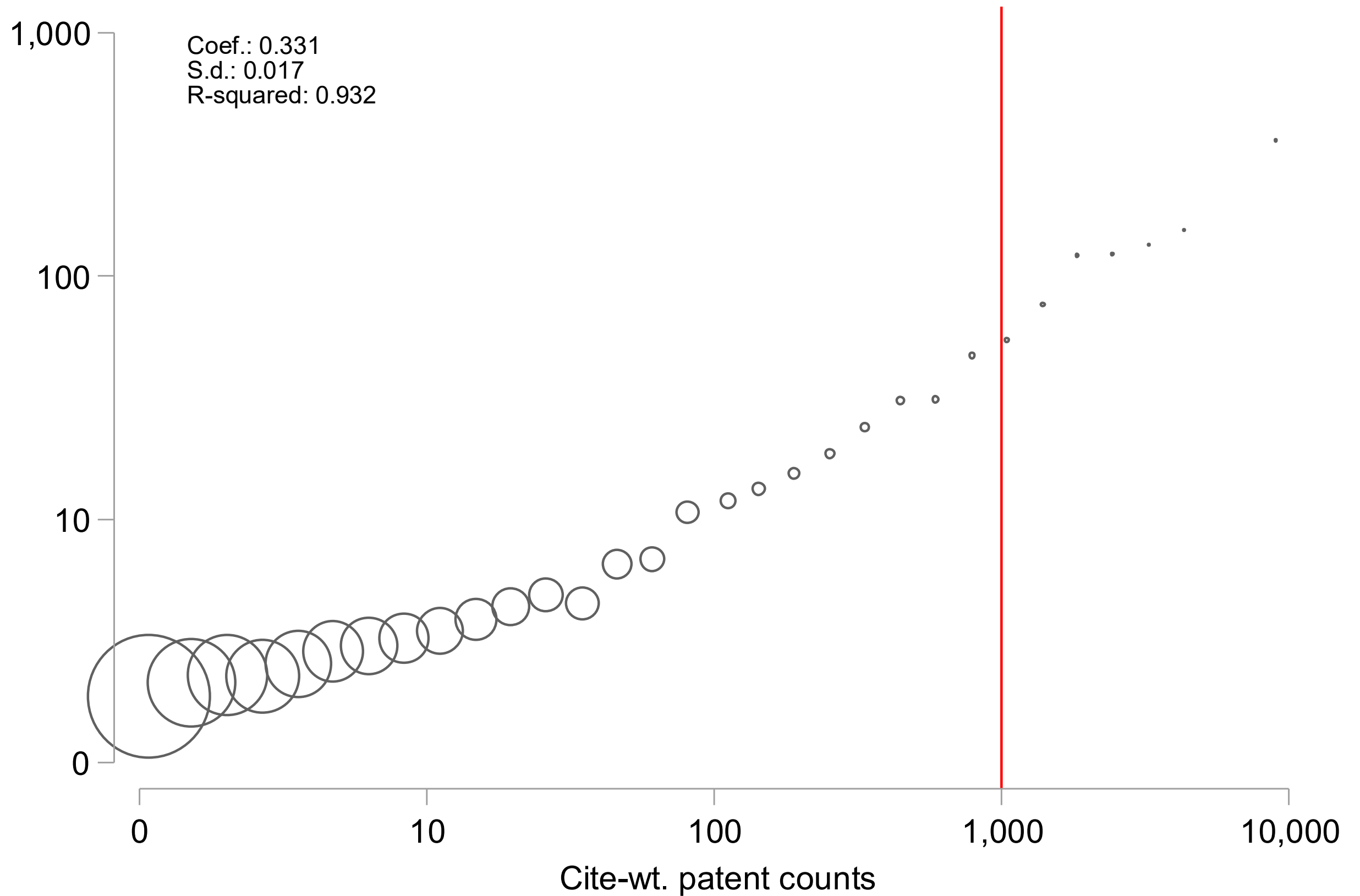


## Panel D: Number of job postings vs employment in OES (by year)



**Notes:** The figure presents descriptions of the changing composition of Burning Glass (BG) job postings over time. Panel A shows the overall number of job postings in BG (left y-axis) against the overall number of job postings in JOLTS (right y-axis). Panel B shows the (approximate) share of job postings in BG requiring a college education by year. Panel C plots a bin-scatter of number for job postings in BG against the number of employed people in Occupational Employment and Wage Statistics (OEWS) program by 6-digit occupation (SOC) and year. Panel D plots the same picture as in Panel A with a bin scatter by year.

**Appendix Figure 10 – Number of normalized citation-weighted patents and job postings associated with novel bigrams**



**Notes:** The figure plots a binned scatterplot of the number of job postings associated with technologies against the number of normalized citation-weighted patents that mention a novel technology. In preparing the figure, we include all technological bigrams with more than 100 citation-weighted patents, as well as a random sample of 1,000,000 novel bigrams with cite-weighted counts between 0 and 99. The size of circles is proportional to the number of bigrams in any given bin. (Unlike in the other figures and regressions, here we do not cap the weights to highlight the concentration of bigrams on the low end of these measures.) The figure presents a binned scatter plot along with the coefficient, standard deviation, and R-squared from a regression of the log of job postings on the log of cite-weighted patent counts corresponding to each technology.