



ECONOMIC RESEARCH
FEDERAL RESERVE BANK OF ST. LOUIS
WORKING PAPER SERIES

Diffusion of new technologies

Authors	Nicholas Bloom Marcela Carvalho, Tarek Hassan, Aakash Kalyani, Josh Lerner, and Ahmed Tahoun
Working Paper Number	2024-009B
Revision Date	September 2024
Citable Link	https://doi.org/10.20955/wp.2024.009
Suggested Citation	Carvalho, N.B.M., Hassan, T., Kalyani, A., Lerner, J., Tahoun, A., 2024; Diffusion of new technologies, Federal Reserve Bank of St. Louis Working Paper 2024-009. URL https://doi.org/10.20955/wp.2024.009

Federal Reserve Bank of St. Louis, Research Division, P.O. Box 442, St. Louis, MO 63166

The views expressed in this paper are those of the author(s) and do not necessarily reflect the views of the Federal Reserve System, the Board of Governors, or the regional Federal Reserve Banks. Federal Reserve Bank of St. Louis Working Papers are preliminary materials circulated to stimulate discussion and critical comment.

The Diffusion of New Technologies

Aakash Kalyani, Nicholas Bloom, Marcela Carvalho, Tarek Hassan,

Josh Lerner, and Ahmed Tahoun¹

August 26, 2024

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 economically impactful new technologies is geographically highly concentrated, more so even than overall patenting: 56% of the most economically 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.

Keywords: Employment, Geography, Innovation, R&D

JEL Classification: O31, O32

The dataset constructed as part of this paper is available at www.techdiffusion.net.

¹ Federal Reserve Bank of St. Louis; Stanford University; Harvard University; Boston University; Harvard University; and London Business School. Bloom, Hassan, and Lerner are affiliated with the National Bureau of Economic Research. The views expressed herein are solely those of the authors and do not necessarily reflect those of the Federal Reserve Bank of St. Louis or the Federal Reserve System. We thank audiences at the Applied Machine Learning webinar, Atlanta Fed, Auburn University, Babson College, the Bank for International Settlements, Baruch College, Bocconi University, CKGSB, College de France, Columbia University, Dartmouth University, Duke University, Durham University, ETH, the Federal Reserve Board, 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 San Diego, 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 2021 NBER Summer Institute, the Fall 2021 NBER EFG meeting, the 2022 Society for Economic Dynamics, and the 2024 American Economic Association annual meetings for helpful comments. Special thanks go to Lisa Kahn for sharing data, Bledi Taska for help on BGT data queries, Gaétan de Rassenfosse, Shane Greenstein, Ben Jones, and Chad Syverson for excellent discussions, and Peter Donets, 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.

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.² Substantial concerns remain, however, 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?).³

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

² 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).

³ 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 new 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 are frequently referenced in these important conversations between firm executives and investors. The most frequently mentioned technologies include “cloud computing,” “smart phone,” and “machine learning.”

Strikingly, the right tail of technologies with the most earnings call mentions also account for the lion’s share of the variation in our job-postings and patenting data. For example, the 276 technologies with more than 100 mentions in earnings calls (the top 22 percent) also appear 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) with application years between 1976 and 2014. In this sense, the innovations that feature prominently in managers’ discussions also have the largest impact on patents and job postings. We therefore pay special attention to these most economically impactful 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 the most economically impactful technologies: 33.3% of patents mentioning any new technology and 42.1% of patents mentioning a new technology with more than 100 mentions in earnings calls 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 the most economically impactful new technologies. Collectively, 56.3% of these most impactful technologies come from just two U.S. locations, Silicon Valley and the Northeast Corridor.⁴ (Locations in California collectively host a remarkable 41.0% of pioneer locations of these most impactful new technologies.) This extreme concentration is particularly important because new technologies alter the composition of the local job postings in their pioneer locations for several decades, 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 U.S., falls by 18.5% in the first decade after its emergence. Nevertheless, the implied years to full dispersion across CBSAs is more than 50 years, well beyond the horizon of most 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 ACS⁵ 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.

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.

⁴ We define the San Jose-Sunnyvale-Santa Clara and San Francisco-Oakland-Hayward CBSAs as Silicon Valley and the Northeast Corridor as New York-Newark-Jersey City, Boston-Cambridge-Newton, Washington-Arlington-Alexandria, and Philadelphia-Camden-Wilmington.

⁵ This share is rising over time, so for example in the 2021 ACS this proportion is 36.4%.

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 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, 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).⁶ 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.⁷ 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.^{8 9}

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. 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

⁶ 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.

⁷ 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).

⁸ Notably, Goldin and Katz (1998) show that the introduction of new manufacturing processes during the early 19th century increased the demand for skilled labor. Krueger (1993) shows that workers who use computers at work earn higher wages.

⁹ Van Reenen (1996) and Kline et al. (2019) study how rents from innovation are shared with employees. We relate to these papers by showing evidence that the economic opportunities stemming from the development of new technologies distribute highly unevenly across space, as opposed to across different actors within a given firm.

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.¹⁰ Important papers by Autor et al. (2024) and Kogan et al. (2024) intersect information from patents, Census job titles, and task descriptions to measure complementarities between innovations and jobs. 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.¹¹

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 and the associated mechanisms. We present our analysis of the skill-broadening results in Section

¹⁰ A notable exception is the work by Abis and Veldkamp (2020), who use job descriptions to identify financial analysis positions that leverage machine learning.

¹¹ Another related strand of work attempts to use large language models to understand the rate and direction of technological change. Recent examples include the work on Kogan et al. (2024) on the impact of technology on labor and Asirvatham's (2024) examination of the speed of technological diffusion.

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.¹² 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.¹³ 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.¹⁴ Because of the importance of the U.S. market, inventors worldwide typically file important discoveries with the USPTO.¹⁵

¹² 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), an abstract idea, a law of nature, nor a 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.

¹³ 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.

¹⁴ This requirement is stipulated in the legal concept of “reduction to practice,” 35 U.S.C. 112(a).

¹⁵ 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 any 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).

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 pre-treat this text in the same way as the patent text, eliminating stop words and extracting bigrams. We then remove any bigram appearing in the *Corpus* (for instance, “equipment used”) from our list of bigrams obtained from patents, leaving us with 1.5 million exclusively “novel” bigrams.¹⁶

¹⁶ At the same time, if the individual words appear in the Corpus, but not in conjunction with each other (e.g., “artificial” and “intelligence” separately, but not as a bigram), we do not delete the phrase.

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.^{17, 18} 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” or “Parkinson’s disease.” Yet others refer to areas that may have seen substantial new developments or management attention but are not new technologies, such as “account management” and “performance metrics.” (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

¹⁷ 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 received as of 2018 by a patent filed in four-digit Combined Patent Classification (CPC) technology class τ in year t . $Avg_{\tau,t}(Citations_{p'})$ is the average number of citations received by all patents filed in technology class τ in year t .

¹⁸ For computational reasons, it is necessary to limit the analysis 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.¹⁹ 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.²⁰ 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

¹⁹ To check the accuracy of this procedure, we conducted a formal human audit following the methodology in Baker et al. (2016). To this end, we developed a detailed coding guide to train three research assistants on the definition of new technologies given above. We asked them to each manually classify a random sample of Wikipedia pages matched to one of our 35,563 novel and influential bigrams from patents. Collectively, the research assistants coded 700 entries. The research assistant’s coding of bigrams that were confidently technological (with a confidence greater or equal to three out of five) corresponded to the answer of the Wikipedia filter 73% of the time. In addition to this human audit, we run robustness tests using a manually reviewed sample of technologies in Section 7.

²⁰ 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.

show our results are robust to using a range of other plausible approaches to defining each bigram’s 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 three CBSAs accounted for 35%, 25%, and 8% of the patents containing a bigram in this period, the first two 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 listed 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).²¹ To gauge the extent to which each technology features in the conversations at these 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

²¹ Some examples of mentions of bigrams in earnings calls are shown in Appendix Table 2.

technological innovation over the past decades. Appendix Table 3 lists all new technologies that 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 impactful technological innovations constructed to date.

Table 2 provides examples of the pioneer locations for several 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).²²

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 and a North American Industry Classification (NAICS) code.²³ 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 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.

²² 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 also 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.”

²³ 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.

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 Glass is missing data for 2008-09, though including the 2007 data has little impact on our results.

Our analysis of job postings thus focuses on the diffusion of technologies with emergence years post-1976 in job postings in the 2010s. That is, technologies with an emergence year of 1980 are thirty years old by the time we see them diffusing in job postings, whereas technologies with an emergence year of 2005 are five years old, and so on. For this reason, we are careful to highlight any differences in the variation across technologies vs. within technologies over time in our analysis below.

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 the respondents who report their occupation in that six-digit SOC code.²⁴ *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.²⁵ *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_τ is a given technology bigram τ 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

²⁴ 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.

²⁵ NAICS codes typically have six nested levels; the four-digit level is referred to as “industry group.”

the technology in the initial firm description or ending boilerplate language, 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 7 for details). As 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 7) requires the worker to use a touch screen to enter data. 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 (4% 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 8 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 6 for an example).²⁶

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

²⁶ Comfortingly, the share of job postings within a given occupation that mentions a new technology with an emergence year post-1979 correlates closely with the share of new job titles created within that occupation since 1980, as identified by Autor et al. (2023) and shown in Appendix Figure 11.

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.²⁷

e. Technologies and earnings calls

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 relatively 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 ($EC \geq 100$) account for about 39 million job postings (or about 77 percent of all job postings mentioning any new technology). On average, each of these technologies is mentioned in 141,634 job postings and 7,682 patents.²⁸

For ease of reference, we sometimes refer to this highly prolific group of new technologies as “economically impactful” new technologies, in the sense that these new technologies take a significant amount of airtime in earnings calls, and feature prominently in both job postings and patents. It includes all the examples from Table 1 (e.g., smartphone, machine learning, hybrid vehicles).

The figure also shows examples of other, less influential, technologies. Those with between 10 and 99 mentions in earnings calls include the pulse oximeter and the liquid chromatograph. On average, each such technology is associated with 26,128 job postings and 4,287 patents. The group

²⁷ Tambe (2014) shows that firms that based in regions with considerable number of workers trained in big data skills experience faster productivity growth, an effect that diminishes as these technologies mature. Below, we show the generalizability of this dissipation result and its slow pace.

²⁸ Interestingly, there is also a clear positive relationship between the numbers of industries in which a given new technology is mentioned and overall earnings call mentions, suggesting that more impactful technologies also tend to be more “general purpose,” in the sense that they are relevant for multiple industries.

with under ten earnings call mentions includes the ultrasonic horn, suction filtration, and NMOS transistors, with on average 1,165 job postings and 3,005 patents.^{29,30}

3. Spatial Concentration of New 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. 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 significantly more concentrated than the distribution of college graduates (22.5%) and the overall workforce (18.9%), but also similar to the concentration of overall patenting activity in the United States (32.4%).³¹

Second, this concentration increases significantly as we condition on increasingly economically impactful technologies as proxied by mentions in earnings calls. Panel A in Table 3 shows that the share of the top-5 CBSAs in patents mentioning new technologies with more than 100 mentions in earnings calls is 42.1%. These prolific CBSAs are San Jose, San Francisco, New York, Seattle, and Boston.

Figure 3 shows the share accounted for by these five prolific CBSAs increases monotonically from 24.6% of patenting relating to relatively low-impact technologies (those mentioned in zero or one earnings calls) to 46.5% of the highest-impact group (new technologies mentioned in 500+ earnings calls). In short, *the most commercially impactful innovations also have the most geographically concentrated origins*.³²

²⁹ 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.

³⁰ Appendix Figure 4 shows the average number of job postings for each category of technology.

³¹ These totals are each for the five CBSAs highest on that individual measure. Only one of the largest CBSAs for patents – New York – is on the top five list for employment, highlighting how population size is not the primary correlate of patenting share.

³² Appendix Table 5 also reports the coefficients for analyses using the top five, three, and one CBSA(s), as well as similar analysis partitioning technologies by the number of associated job postings.

Interestingly, to preview our results below, this extreme concentration of economically impactful innovation is the only significant difference that we document between more and less economically impactful innovations. Aside from their concentrated origins, less impactful technologies appear to evolve and spread similarly to their more impactful counterparts.

In the same vein, Figure 4 shows the distribution of pioneer CBSAs – the urban areas that account for a majority of *early* patenting of economically impactful technologies (again, those with more than 100 earnings call mentions, $EC \geq 100$). Panel A of Figure 4 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 and San Francisco 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, Boston, Washington DC, and Philadelphia, 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 high concentration of the most economically impactful innovative activity within America over the last decades.

These pioneer locations tend to have highly educated workforces and a high density of university activity. For each CBSA- technology pair (e.g., “smart 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 all cases, there is a strong association between measures of education/university presence and per capita patents relating to new technologies. Interestingly, these associations are significantly more pronounced when we condition on economically impactful new technologies. 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 key new technologies.³³

We show evidence below that this concentration of innovation in a handful of urban centers engenders large and persistent regional disparities in economic opportunity, as measured by job

³³ 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.)

postings in local labor markets. In this sense, a handful of U.S. locations appear to have a comparative advantage in developing technologies that most impact 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 τ 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 τ 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 . $Normalized\ share_{c,\tau,t}$, therefore, measures the regional over- or under-representation 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.³⁴

In Figure 5, we present a series of maps displaying the spread of job postings mentioning economically impactful new technologies ($EC \geq 100$). The blue circles identify the same pioneer locations as in Figure 4, but now superimpose purple dots that show the intensity of the normalized share of job postings relating to these new technologies 0-5, 6-10, 11-20, and 21-30 years after the 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 new technologies gradually spread across space (region broadening). Second, there is a remarkable alignment between the CBSAs that pioneer early development in 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

³⁴ Throughout, we cap this variable at the 99th percentile of non-zero observations.

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 τ .³⁵ 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 new technologies is 4.69, which suggests that technology job postings relating to these new 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:

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

where $CV_{\tau,t}$ is the coefficient of variation across CBSAs for technology bigram τ in year t , and $(t - t_{0,\tau})$ is the number years since emergence of technology bigram τ in year $t_{0,\tau}$ (capped at 30 years, given we have little data for technologies older than 30 years). δ_{τ} denotes a full set of technology bigram fixed effects, which we constrain to sum to zero, so that the intercept α_0 measures the average coefficient of variation in the year of emergence.³⁶ The slope coefficient, β_{RB} , measures the speed of decay of this concentration with each passing year since emergence. Panel A, column 1 reports estimates for economically impactful technologies ($EC \geq 100$). 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).³⁷

³⁵ Appendix Table 8 summarizes the data used in this and subsequent regression analyses.

³⁶ 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, we exclude technology bigrams with less than 1,000 job postings. In practice, these adjustments have little impact on our estimates (see Section 7).

³⁷ 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). However, note that the dependent variable is already normalized to account for any time trends in the overall coverage of job postings: By construction, the coefficient of variation of the overall job postings in our database is 0 for all t , meaning that our measure of the diffusion of technology job postings is immune to variations over time in the share of jobs covered by BG or the shares of regional labor markets covered.

In column 1, we find that on average, in the year of emergence, the coefficient of variation is 5.58 (significantly greater than 0). With each additional year since emergence, this coefficient of variation decreases by 0.068 (s.e.=0.026) points (or 1.22%). Taking these estimates at face value suggests that technology job postings on average take 82 years to fully disperse across the U.S. (This latter projection is of course considerably out of sample.)

Panel A of Figure 6 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.³⁸ This is corroborated by Panel B of Figure 6, which shows the average of the coefficient of variation separately for five age groups (3-10, 11-15, 16-20, 21-25, and 26-30+ years since emergence), also reporting standard error bands for these estimates. Again, we see a clear decline, with no clear acceleration or deceleration in the rate of decline.

In column 2 of Panel A of Table 4, we show this pattern is almost identical when we include all new technologies in our sample. Appendix Table 19 tests explicitly for differences in the rate of spread between technologies with fewer and greater than 100 earnings call mentions, finding no economically significant differences.

In column 3 of Panel A 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 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 one (similar to the

³⁸ In Appendix Table 9, Panels A and B add a quadratic term to all of the specifications documenting Region Broadening in Table 4. The coefficient on the quadratic term is statistically indistinguishable from zero in five of the six specifications. The only exception is Column 3 of Panel A, where the coefficient on the quadratic term is significant at the 10% level (0.002, s.e.=0.001). Panels C and D expands on this theme by adding an interaction with a dummy for technologies older than 20 years to all the specifications shown in Table 4. Of the six specifications shown, five show no statistically significant difference in the marginal effects of a technology's age (year since emergence) for older and younger technologies.

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 relating to economically impactful technologies ($EC \geq 100$); columns 3 and 4 show results for all technologies.

In column 1, we see that pioneer locations enjoy a significant pioneer advantage on average: 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 3.0% (0.032/1.084). The initial advantage of the pioneering locations for job postings relating to the economically impactful technologies they develop thus lasts for decades, with an implied 34 years to zero advantage.³⁹

In column 3, we include all technologies and again find an almost identical pattern -- albeit with a somewhat larger point estimate for the pioneer advantage in 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, $Pioneer\ Neighbor_{c,\tau}$, 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

³⁹ In Appendix Table 10, 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.

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 new 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 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.^{40, 41}

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

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

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

⁴⁰ As an example, Appendix Table 11 shows the list of top occupations by share of job postings for some of our top technologies (see Section 2 of the Data Appendix for details).

⁴¹ The BG data also includes an indicator for a college requirement for a subset of observations. However, since this subset is quite limited, we prefer using SOC codes to generate this variable.

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.⁴²

We find that the geographic concentration of low-skill jobs (in column 2) decreases 1.1 percentage points or 41% ($0.038/0.027 - 1$) faster than that of high-skill jobs (in column 1). This difference is statistically significant at the 1% level, as we report in the label of Table 6 (and in the labels of subsequent analyses where we can compare coefficients across equations). Figure 7, 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 2.2 percentage points per year, the latter erodes at a faster 3.2 percentage points. These estimates imply it takes 45 years for a pioneer location's advantage in high-skill jobs to erode, whereas that for low-skill jobs lasts only 31 years.⁴³

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.

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

⁴² Results are almost identical when using a tripartite division of skill levels, as Appendix Table 12 shows.

⁴³ Both results are again almost identical when we repeat these estimations for the subset of technologies with $EC \geq 100$ in Appendix Table 20.

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 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 are typically used in similar context to these keywords – in effect, using the embedding model like a custom-trained thesaurus.⁴⁴ 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 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 13 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 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 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

⁴⁴ 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.

large differences: technology-using job postings spread out 157.1% ($=0.036/0.014-1$) faster than postings that involve technology RDP jobs. This difference is again 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 5 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 158.1% more such jobs in the pioneer location in the year of emergence) and dissipates significantly over time (-0.030 , $s.e.=0.012$). By contrast, RDP job postings are more concentrated in pioneer locations initially (197.4% 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 7, 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 Masters- and PhD-level 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 τ 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 τ in year t all associated job postings are in an occupation o , then its skill level is 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 τ 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 (β_{SB}) denotes this skill level's average speed of decay with each passing year since emergence. Column 1 of Panel A reports results for economically impactful 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 in 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.

Panel A of Figure 8 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.

Panel B shows the average of the skill level separately for five age groups (again, 3-10, 11-15, 16-20, 21-25, and 26-30+ years since emergence), as well as standard error bands for these estimates. Visual inspection of the non-parametric specifications in Figure 8 again suggests a linear

relationship. Column 2 of Table 7 shows almost identical results for the broader sample with all technologies.

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.⁴⁵ To address this concern, column 3 shows that 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 14 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). 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.

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.,

⁴⁵ 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.

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 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 or experience – over time, less educated workers may acquire training or experience 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”). 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 15).

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. At the same time, training / experience requirements with the new technology 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 (Appendix Figure 8).

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, β_{SB} . Column 1 reproduces our estimate from column 2 of Panel A, Table 7 for comparison (-0.288, s.e.=0.079). Column 4 shows that controlling for the inverse hyperbolic sine of the share of RDP jobs in the same technology attenuates this estimate by about 22% to -0.224 (s.e.=0.055). 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.⁴⁶ 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 (in fact, an increase) of our estimate of β_{SB} . In this sense, changes in training and experience in the technology cannot account for the pattern of skill broadening observed in the data.

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

Finally, 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 new technologies spread across occupations, firms, and industries, we extend the definition of *Normalized share*_{*c,t,t*} to NAICS four-digit industries, SOC six-digit

⁴⁶ 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.

occupations, and firms for each technology (τ) and time (t), calculating the normalized share of job postings in each industry, occupation, and firm that mention a given new technology.⁴⁷ We then measure the coefficient of variation of *Normalized share* _{c,τ,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.⁴⁸ This step focuses attention on 10,496 larger firms, which on average post 865 job postings 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 (τ) and time (t) on the year since emergence. Column 4 shows our already established results for locations for comparison. We find that while there is a decline in concentration as measured by the 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). 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.⁴⁹ 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.⁵⁰ In 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 in the year of emergence for pioneer firms. Over time,

⁴⁷ 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.

⁴⁸ 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.

⁴⁹ The decay rates across CBSAs are 0.016 (0.004), 0.011 (0.003), and 0.003 (0.002) higher than industries, occupations, and firms, respectively. These coefficients are statistically significant at the 1%, 1%, and 20% level, respectively.

⁵⁰ 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 – and Appendix Table 16 for some examples.

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 impactful 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.⁵¹

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

⁵¹ 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.

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 17).

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 technology already identified in our bigram-based analysis. Appendix Table 18 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”). 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 material 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 economically impactful technologies dropped under this audit (altogether 63 of 276 technologies with $EC \geq 100$). 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 (i) with 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. We also re-run our entire analysis after collapsing technology bigrams at the technology (Wikipedia title) level. Again, none of these variations materially affect our results.

In the final line of the panel, we re-calculate our list of economically impactful technologies using an emergence-year-normalized number of earnings calls mentions: For each technology bigram i with year of emergence t_0 , we divide the number of earnings calls appearances by the average number of earnings calls for all bigrams with year of emergence t_0 . This adjustment alters our list of influential technologies by controlling for the different number of years that the various bigrams had to be mentioned in earnings calls. Again, doing so has little influence on 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 14 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 4) 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 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 influential 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 most economically impactful innovations, 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.
- Asirvatham, Hemanth. "The dynamo is not the computer: A historical exploration of what makes tech adoption faster." Senior thesis, Harvard College (2024).
- 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., Caroline Chin, Anna Salomons, and Bryan Seegmiller. "New frontiers: the origins and content of new work, 1940-2018." *Quarterly Journal of Economics*, qjae008 (2024).

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.

- 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 (2008).

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.

Kline, Patrick, Neviana Petkova, Heidi Williams, and Owen Zidar. "Who profits from patents? Rent-sharing at innovative firms." *Quarterly Journal of Economics* 134 (2019): 1343-1404.

Kogan, Leonid, Dimitris Papanikolaou, Lawrence D. Schmidt, and Bryan Seegmiller. "Technology 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. "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) (1968): 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).
- Van Reenen, John. "The creation and capture of rents: Wages and innovation in a panel of U. K. companies." *Quarterly Journal of Economics* 111 (1996): 195-226.
- 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 of emergence

Emergence year	Wikipedia title (technology)	Technology bigrams	Number of job postings
1979	Hard disk drive	hard disk; disk drive	34,211
1980	Barcode reader	barcode reader; code reader; code scanner; barcode scanner	43,279
1981	Laser diode	emitting laser; diode laser; semiconductor laser; laser diode	7,284
1982	Personal computer	personal computer	1,752,726
1983	Flat-panel display	panel display; flat panel	27,369
1984	User interface	user interface	747,586
1985	Mobile phone	mobile telephone; cellular telephone; phones mobile; cellular phone; mobile phone; cell phone	1,832,787
1986	Facial recognition system	frt system; recognition software; recognition system; recognition technology; facial recognition	25,109
1987	Digital video	digital video	88,887
1988	Model organism	animal model	24,722
1989	Mobile device	held computer; computer device; handheld computer; mobile device	1,046,079
1990	Debit card	cards debit; card debit; debit card	260,282
1991	Flash memory	flash device; nand flash; flash memory	22,882
1992	Machine learning	learning algorithm; machine learning	491,252
1993	Financial instrument	financial instrument	43,944
1994	Active users	active user	39,671
1995	Hybrid electric vehicle	hybrid electric	8,207
1996	Digital content	digital content	144,775
1997	Multicore processor	multi core; core processor	29,643
1998	Information privacy	data protection	176,110
1999	Unmanned aerial vehicle	aerial vehicle; unmanned aerial	24,148
2000	Transaction account	transaction account	13,012
2001	Smartphone	smart phone	910,856
2002	Online game	online game	15,254
2003	Social networking service	networking site; social networking	244,610
2004	Electronic discovery	electronic format	56,438
2005	LED circuit	led driver	2,575
2006	Augmented reality	augmented reality	20,537
2007	Self-driving car	autonomous vehicle	18,641

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). Column 4 lists the number of job postings that the bigram appears in. 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.” Column 4 reports the number of job postings associated with the technology. 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			
Economically impactful	1,044,351	42.1%	San Jose-Sunnyvale-Santa Clara, CA San Francisco-Oakland-Hayward, CA New York-Newark-Jersey City, NY-NJ-PA Seattle-Tacoma-Bellevue, WA Boston-Cambridge-Newton, MA-NH
All New Technologies	1,623,800	33.3%	San Jose-Sunnyvale-Santa Clara, CA San Francisco-Oakland-Hayward, CA New York-Newark-Jersey City, NY-NJ-PA Los Angeles-Long Beach-Anaheim, CA Boston-Cambridge-Newton, MA-NH
All Patents	3,146,114	32.4%	San Jose-Sunnyvale-Santa Clara, CA New York-Newark-Jersey City, NY-NJ-PA San Francisco-Oakland-Hayward, CA Los Angeles-Long Beach-Anaheim, CA Chicago-Naperville-Elgin, IL-IN-WI
Most Cited	1,044,351	32.7%	San Jose-Sunnyvale-Santa Clara, CA San Francisco-Oakland-Hayward, CA New York-Newark-Jersey City, NY-NJ-PA 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%	New York-Newark-Jersey City, NY-NJ-PA Los Angeles-Long Beach-Anaheim, CA Chicago-Naperville-Elgin, IL-IN-WI Washington-Arlington-Alexandria, DC-VA-MD-WV San Francisco-Oakland-Hayward, CA
Employed	156.5 million	18.9%	New York-Newark-Jersey City, NY-NJ-PA 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 an economically impactful/new technology patent if it mentions at least one bigram associated with an economically impactful/new technology 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 an economically impactful technology. CBSAs in **bold** are those in the top five for patents which mention economically impactful technologies.

Table 4 – Region broadening

Panel A: Main specifications			
	<i>Coefficient of Variation_{τ,t}</i>		
Sample	EC \geq 100	All	
	(1)	(2)	(3)
<i>Years since emergence_{τ,t}</i>	-0.068*** (0.026)	-0.065*** (0.024)	-0.153*** (0.012)
Constant (CV at $t=t_{\tau,0}$)	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
Years to zero CV	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_{τ,t}</i>	-1.724*** (0.136)	-1.117*** (0.088)	-173.965** (77.247)
Cons (concentration at $t=t_{\tau,0}$)	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
Years to zero CV	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 one (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 of job postings mentioning that technology in that year, capped at 100. The normalized share of job postings is capped at the 99th 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_{c,t,t}</i>			
	<i>EC_τ ≥ 100</i>		All	
	(1)	(2)	(3)	(4)
<i>Pioneer_{c,τ}</i>	0.311*** (0.076)	1.084*** (0.309)	1.321*** (0.254)	1.282*** (0.243)
<i>Pioneer_{c,τ} * Years since emg_{τ,t}</i>		-0.032** (0.013)	-0.035*** (0.011)	-0.034*** (0.010)
<i>Pioneer Neighbor_{c,τ}</i>				0.158*** (0.057)
<i>Pioneer Neighbor_{c,τ} * Years since emg_{τ,t}</i>				-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)
Implied years to zero advantage		33.88	37.74	38.26

Notes: This table reports results from regressions of the *Normalized Share_{c,t,t}* (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_{c,τ}* 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 of job postings mentioning that technology in that year, capped at 100. The normalized share of job postings is capped at the 99th 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 β_P is the coefficient on *Pioneer_{c,τ}* and β_D is the coefficient of *Pioneer_{c,τ} * Years since emg_{τ,t}*.

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				
Sample:	$\log(\text{Coefficient of Variation})_{\tau,t}$			
	All			
	(1)	(2)	(3)	(4)
	High-Skill Job Postings	Low-Skill Job Postings	RDP Job Postings	Use Job Postings
<i>Years since emergence</i> $_{\tau,t}$	-0.027*** (0.002)	-0.038*** (0.003)	-0.014*** (0.002)	-0.036*** (0.002)
R-squared	0.837	0.845	0.736	0.883
N	8,069	8,069	6,033	6,033
Bigram FE	YES	YES	YES	YES
Std. Errors (cluster)	Wiki Title	Wiki Title	Wiki Title	Wiki Title
Panel B: Pioneer advantage regressions				
Sample:	$\text{Normalized Share}_{c,\tau,t}$			
	All			
	(1)	(2)	(3)	(4)
	High-Skill Job Postings	Low-Skill Job Postings	RDP Job Postings	Use Job Postings
<i>Pioneer</i> $_{c,\tau}$	1.319*** (0.233)	1.127*** (0.255)	1.974*** (0.595)	1.581*** (0.301)
<i>Pioneer</i> $_{c,\tau} * \text{Years since emg}$ $_{\tau,t}$	-0.029*** (0.010)	-0.036*** (0.010)	-0.021 (0.027)	-0.030** (0.012)
R-squared	0.016	0.012	0.003	0.020
N	8,581,946	8,395,144	5,618,723	7,769,596
Bigram FE	YES	YES	814	837
CBSA FE	YES	YES	YES	YES
Year FE	YES	YES	YES	YES
Std. Errors (cluster)	Wiki Title	Wiki Title	YES	YES
<i>Rate of decline per year</i>	-0.022 0.004	-0.032 0.003	-0.011 0.011	-0.019 0.005
Implied years to zero advantage	45.48	31.31	104.28	51.76

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 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 of 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 -0.011 (S.E. = 0.003, p-val. = 0.001). The difference between the coefficient on *Years since emergence* $_{\tau,t}$ in columns 3 and 4 is -0.023 (S.E. = 0.003, p-val = 0.000). Both differences are thus statistically distinguishable from zero. Panel B reports results from regressions of the $\text{Normalized share}_{c,\tau,t}$ (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 99th 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 0.010 (S.E. = 0.005, p-val = 0.026). Similarly, the difference between the rate of decline per year in columns 3 and 4 is 0.007 (S.E. = 0.007, p-val = 0.518). The rate of decline per year is calculated as $\frac{\beta_D}{\beta_P}$, where β_P is the coefficient on *Pioneer* $_{c,\tau}$ and β_D is the coefficient of *Pioneer* $_{c,\tau} * \text{Years since emg}$ $_{\tau,t}$.

Table 7 – Skill broadening

Panel A: Main specifications				
<i>Share College Educated_{τ,t} * 100</i>				
Sample	EC >= 100	All		
	(1)	(2)	(3)	(4)
Constant (Sh. Col. Ed. at $t=t_{\tau,0}$)	57.078*** (2.135)	59.095*** (1.794)	57.475*** (2.294)	63.898*** (0.840)
<i>Years since emergence_{τ,t}</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
Implied years to average skill	117.23	100.03	124.37	68.08
Panel B: Alternative measures of skill				
	(1)	(2)	(3)	
	<i>Years of Schooling_{τ,t}</i>	<i>Share Post Graduates_{τ,t} * 100</i>	<i>Average Wage_{τ,t}</i>	
Constant (Skill at $t=t_{\tau,0}$)	15.504*** (0.047)	22.617*** (0.456)	75,521.317*** (840.562)	
<i>Years since emergence_{τ,t}</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	
Implied years to average skill	77.59	78.03	69.72	

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 τ 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 of 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 $-(\text{Constant} - \text{Average Population Skill})/\beta_{SB}(\text{Years since emergence}_{\tau,t})$, 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	Share College Educated $d_{\tau,t}$				
	All				
	(1)	(2)	(3)	(4)	(5)
<i>Years since emergence$_{\tau,t}$</i>	--0.288*** (0.079)	--0.231*** (0.056)	--0.268*** (0.063)	--0.224*** (0.055)	--0.340*** -0.069
<i>Share of R&D Postings$s_{\tau,t}$, IHS</i>		6.938*** (0.366)			
<i>Share of Produce Postings$s_{\tau,t}$, IHS</i>			5.528*** (0.388)		
<i>Share of RDP Postings$s_{\tau,t}$, IHS</i>				6.795*** (0.381)	
<i>Share Training Required$_{\tau,t}$, IHS</i>					5.553*** -0.459
<i>Constant</i>	59.095*** (1.794)	43.482*** (1.499)	46.661*** (1.670)	39.136*** (1.662)	39.619*** (2.413)
R-squared	0.019	0.432	0.276	0.407	0.236
N	8,347	8,347	8,347	8,347	8,347
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 2. Columns 2-5 add additional controls: the inverse hyperbolic sine (IHS) of 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). Observations are weighted by the square root of the total number of job postings mentioning that technology in that year, capped at 100. Standard errors are clustered by Wikipedia title (technology).

Table 9 – Broadening and pioneer advantage across different dimensions

Panel A: Broadening				
	<i>Coefficient of Variation_{τ,t}</i>			
	(1)	(2)	(3)	(4)
	Industries	Occupations	Firms	CBSAs
<i>Years since emergence_{τ,t}</i>	-0.018 (0.017)	-0.056*** (0.015)	-0.354*** (0.038)	-0.153*** (0.012)
Cons (CV at $t=t_{\tau,0}$)	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
<i>Rate of decline per year</i>	-0.004 (0.003)	-0.007 (0.002)	-0.016 (0.001)	-0.018 (0.001)
Years to zero CV	273.38	145.29	62.21	54.07
Panel B: Pioneer advantage				
	<i>Normalized Share_{n,τ,t}</i>			
	(1)	(2)	(3)	
	Industries	Firms	CBSAs	
<i>Pioneer_{n,τ}</i>	6.504*** (1.751)	20.935*** (4.883)	1.321*** (0.254)	
<i>Pioneer_{n,τ} * Years since emg_{τ,t}</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	
<i>Rate of decline per year</i>	-0.013 (0.008)	-0.023 (0.004)	-0.027 (0.003)	
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 the coefficient of variation calculated over *Normalized Share_{n,τ,t}* 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_{n,τ,t}* 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 of 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 0.016 (0.004)***, 0.011 (0.003)***, and 0.003 (0.002) higher than across 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_{n,τ}* is 19.614 (6.604)*** and 5.183 (1.659)*** higher for firms and industries than for CBSAs. For the coefficient on *Pioneer_{i,τ} * Years since emg_{τ,t}*, 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, 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. 3]	Skill Broadening [Table 7, Panel A, col. 4]
Estimate/Coefficient:	Share of Top-5 CBSAs	β_{RB}	$\beta_{RB}^{High\ skill} - \beta_{RB}^{Low\ skill}$	β_D/β_P	β_{SB}
Panel A: Influential patents					
Baseline: At least 1,000 cite-wt. patents	42.1%	-0.153*** (0.012)	-0.011*** (0.003)	-0.027*** (0.003)	-0.493*** (0.036)
At least 1,250 cite-wt. patents	42.2%	-0.150*** (0.012)	-0.010*** (0.003)	-0.027*** (0.003)	-0.488*** (0.037)
At least 1,500 cite-wt. patents	42.4%	-0.146*** (0.012)	-0.011*** (0.003)	-0.027*** (0.003)	-0.500*** (0.038)
At least 1,750 cite-wt. patents	42.5%	-0.146*** (0.013)	-0.009*** (0.003)	-0.027*** (0.004)	-0.501*** (0.040)
At least 2,000 cite-wt. patents	42.5%	-0.147*** (0.013)	-0.009*** (0.003)	-0.027*** (0.004)	-0.494*** (0.040)
Panel B: Phrase Length					
Baseline: Bigrams	42.1%	-0.153*** (0.012)	-0.011*** (0.003)	-0.027*** (0.003)	-0.493*** (0.036)
Bigrams and trigrams	42.0%	-0.157*** (0.011)	-0.011*** (0.003)	-0.025*** (0.004)	-0.494*** (0.035)
Bigrams, trigrams, and unigrams (economically impactful only)	37.0%	-0.124*** (0.008)	-0.012*** (0.003)	-0.029*** (0.003)	-0.486*** (0.033)
Panel C: Human Audit					
Baseline: Bigrams	42.1%	-0.153*** (0.012)	-0.011*** (0.003)	-0.027*** (0.003)	-0.493*** (0.036)
Human-audited bigrams (economically impactful only)	44.4%	-0.142*** (0.019)	-0.014*** (0.005)	-0.025*** (0.009)	-0.574*** (0.061)
Panel D: Alternative emergence years					
Baseline: At least 100 cite-wt. patents	42.1%	-0.153*** (0.012)	-0.011*** (0.003)	-0.027*** (0.003)	-0.493*** (0.036)
At least 1 cite-wt. patent	41.9%	-0.157*** (0.011)	-0.013*** (0.003)	-0.026*** (0.003)	-0.465*** (0.034)
At least 200 cite-wt. patents	42.2%	-0.154*** (0.013)	-0.012*** (0.003)	-0.028*** (0.003)	-0.511*** (0.042)
50% of total cite-wt. patents	39.2%	-0.144*** (0.009)	-0.011*** (0.003)	-0.028*** (0.008)	-0.466*** (0.027)
Panel E: Alternative weighting schemes					
Baseline: min(100, sqrt(# postings))	NA	-0.153*** (0.012)	-0.011*** (0.003)	-0.027*** (0.003)	-0.493*** (0.036)
Unweighted regression	NA	-0.194*** (0.014)	-0.010*** (0.002)	-0.019*** (0.005)	-0.532*** (0.043)

All bigrams with job postings	NA	-0.209*** (0.014)	-0.009*** (0.002)	-0.020*** (0.003)	-0.490*** (0.046)
Log-wt: min(100, log(# postings))	NA	-0.178*** (0.013)	-0.011*** (0.003)	-0.022*** (0.004)	-0.515*** (0.039)
Bigrams collapsed into technologies	NA	-0.156*** (0.010)	-0.013*** (0.002)	-0.026*** (0.005)	-0.490*** (0.033)
100+ normalized EC counts	NA	-0.149*** (0.015)	-0.013*** (0.004)	-0.024*** (0.006)	-0.532*** (0.044)

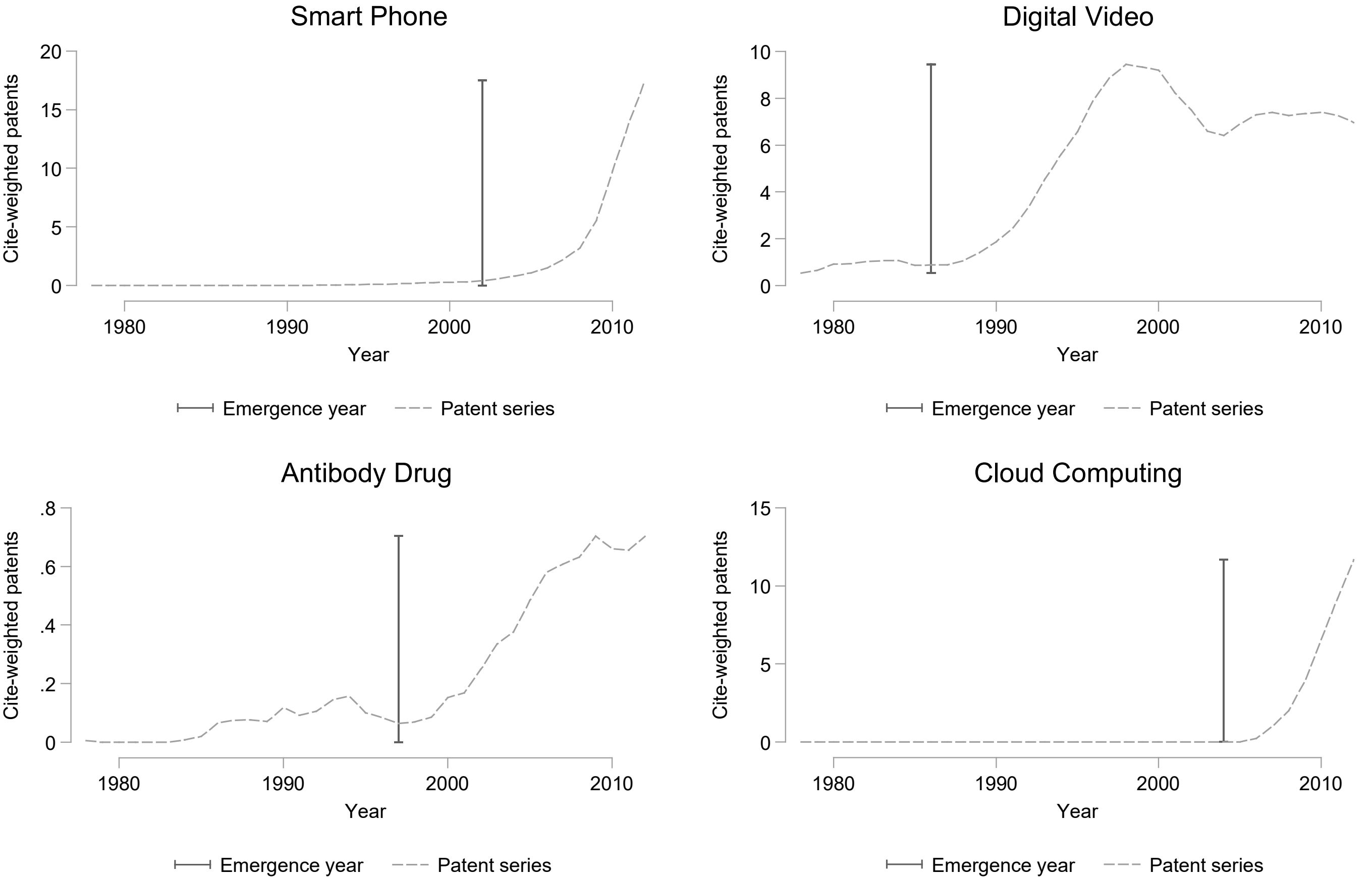
Notes: This table reports robustness checks to our primary (“Baseline”) results. Panel A reports robustness to changing the threshold for defining bigrams associated with “influential innovations.” In the baseline, we retain only those that appear in patents accumulating a total of at least 1,000 weighted 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 (economically impactful only)” further adds unigrams. In Panel C, “Human audited bigrams (economically impactful only)”, we rely on human reading to determine whether or not a bigram describes a technology, instead of the Wikipedia filter, reporting results from including only those technology bigrams that survived the human auditing process. 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 (b) 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 of 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 row “Bigrams collapsed into technologies” replicates our results when all bigrams associated with a given Wikipedia title (technology) are collapsed into the technology. The last row, “100+ normalized EC counts” replicates our results with bigrams that cumulate more than 100 normalized earnings calls mentions. The rows reporting unigrams in Panel B and human-audited bigrams in Panel C report results only using economically impactful technologies. See the original regressions for full details.

Table 11 – Robustness checks: Alternative specifications of standard errors

	Coefficient of Variation (1)	Coefficient of Variation (2)	Normalized Share (3)	Share College Educated (4)
	Region Broadening [Table 4, Panel A, col. 3]	Region Broadening by Skill [Table 6, Panel A, col 1,2]	Pioneer Persistence [Table 5, col. 3]	Skill Broadening [Table 7, Panel A, col. 4]
<i>Years since emergence_{τ,t}</i> (<i>High Skill</i> in col. 2)	-0.153	-0.027		-0.493
[baseline] Cluster, Wikipedia Title level	(0.012) ***	(0.002) ***		(0.036) ***
Cluster, Bigram level	(0.009) ***	(0.002) ***		(0.028) ***
Cluster, Year level	(0.018) ***	(0.004) ***		(0.042) ***
Bootstrap (500 replications)	(0.009) ***	(0.002) ***		(0.027) ***
<i>Years since emergence_{τ,t}</i> (<i>Low Skill</i>)		-0.038		
[baseline] Cluster, Wikipedia Title level		(0.003) ***		
Cluster, Bigram level		(0.002) ***		
Cluster, Year level		(0.005) ***		
Bootstrap (500 replications)		(0.002) ***		
<i>Pioneer_{i,τ}</i>			1.321	
[baseline] Cluster, Wikipedia Title level			(0.254) ***	
Cluster, Bigram level			(0.213) ***	
Cluster, Year level			(0.059) ***	
Cluster, CBSA level			(0.203) ***	
Cluster, State level			(0.190) ***	
Cluster, CBSA-Wikipedia Title levels			(0.277) ***	
Bootstrap (500 replications)			(0.214) ***	
<i>Pioneer_{i,τ} * Years since emergence_{τ,t}</i>			-0.035	
[baseline] Cluster, Wikipedia Title level			(0.011) ***	
Cluster, Bigram level			(0.009) ***	
Cluster, Year level			(0.003) ***	
Cluster, CBSA level			(0.007) ***	
Cluster, State level			(0.006) ***	
Cluster, CBSA-Wikipedia Title levels			(0.011) ***	
Bootstrap (500 replications)			(0.008) ***	
$\beta(Pioneer_{i,\tau} * Years\ since\ emergence_{\tau,t})$ $/\beta(Pioneer_{i,\tau})$			-0.027	
[baseline] Cluster, Wikipedia Title level			(0.003) ***	
Cluster, Bigram level			(0.003) ***	
Cluster, Year level			(0.001) ***	
Cluster, CBSA level			(0.003) ***	
Cluster, State level			(0.003) ***	
Cluster, CBSA-Wikipedia Title levels			(0.004) ***	
Bootstrap (500 replications)			(0.003) ***	
Bigram FE	YES	YES	YES	YES
Skill FE	NA	YES	NA	NA
CBSA FE	NA	NA	YES	NA
Year FE	NA	NA	YES	NA

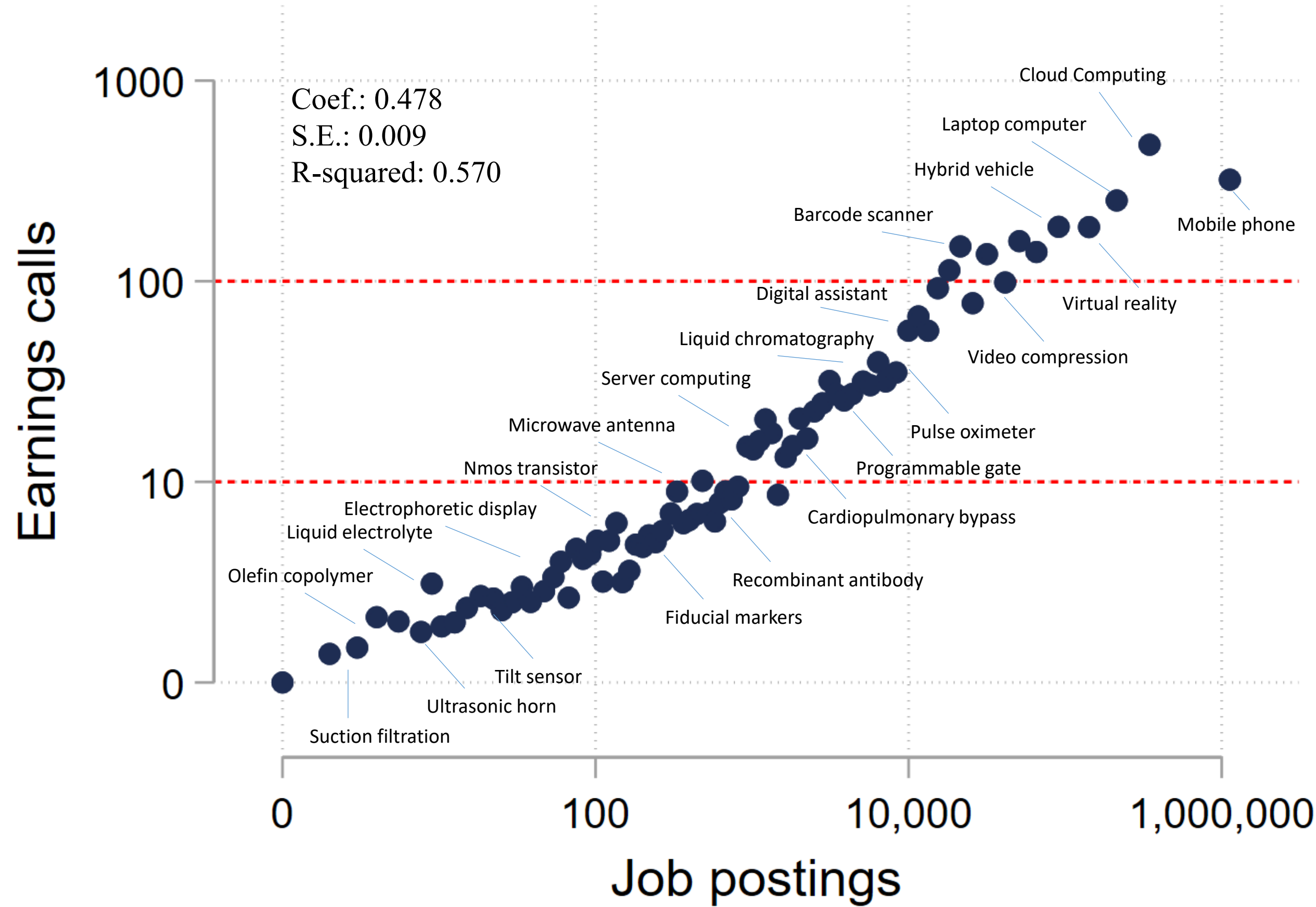
Notes: This table reports results from varying specifications for standard errors corresponding to coefficient estimates for our main results – region broadening, region broadening by skill, pioneer persistence and skill broadening. The statistical significance of coefficients is indicated by the asterisks 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 500 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



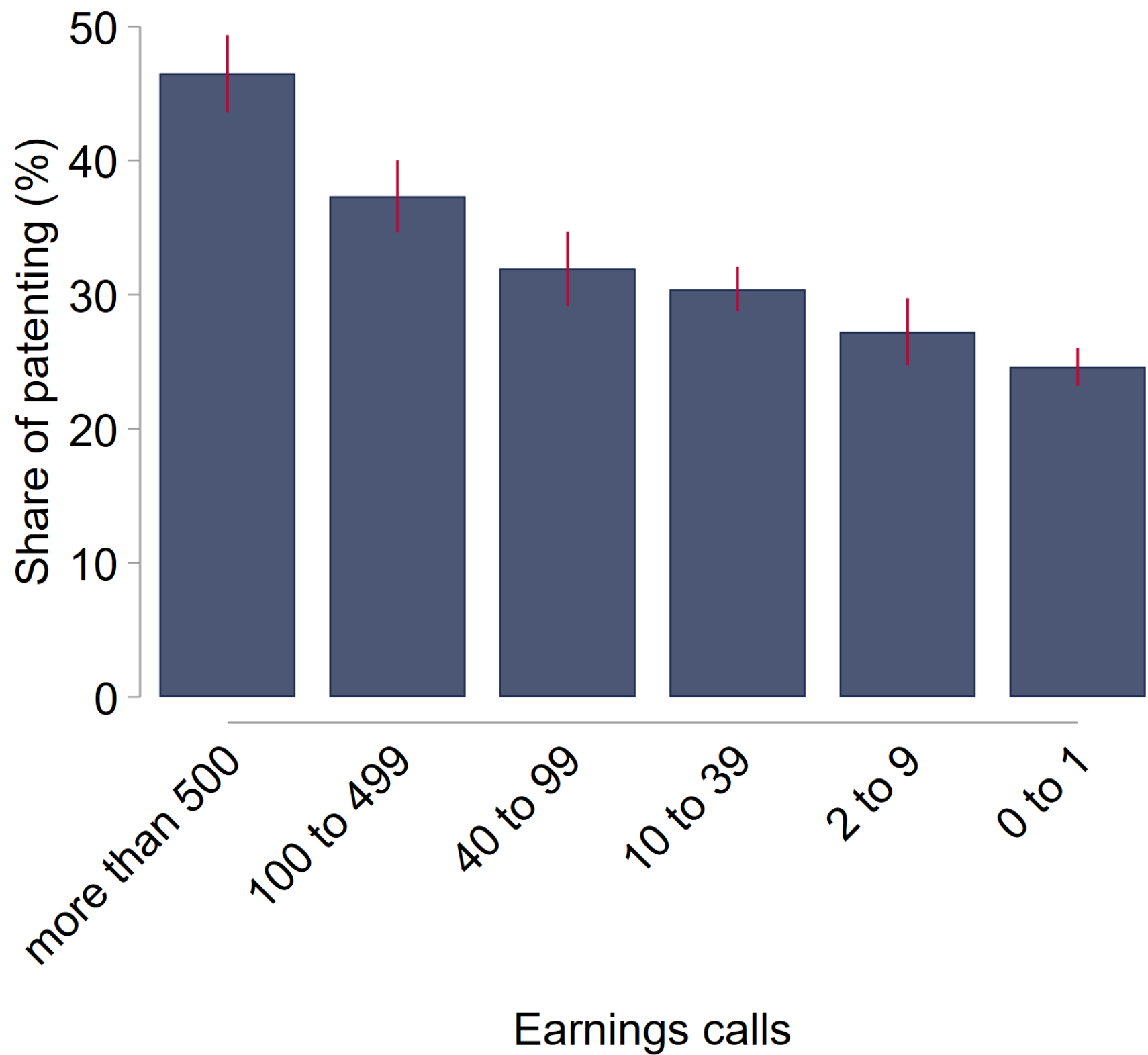
Notes: The figure shows four examples of the attribution of emergence years. In each example, the time series plots the smoothed number of cite-weighted patents associated with the technology by year of application of the patent. For each bigram, we mark the emergence year as the first year in which (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.

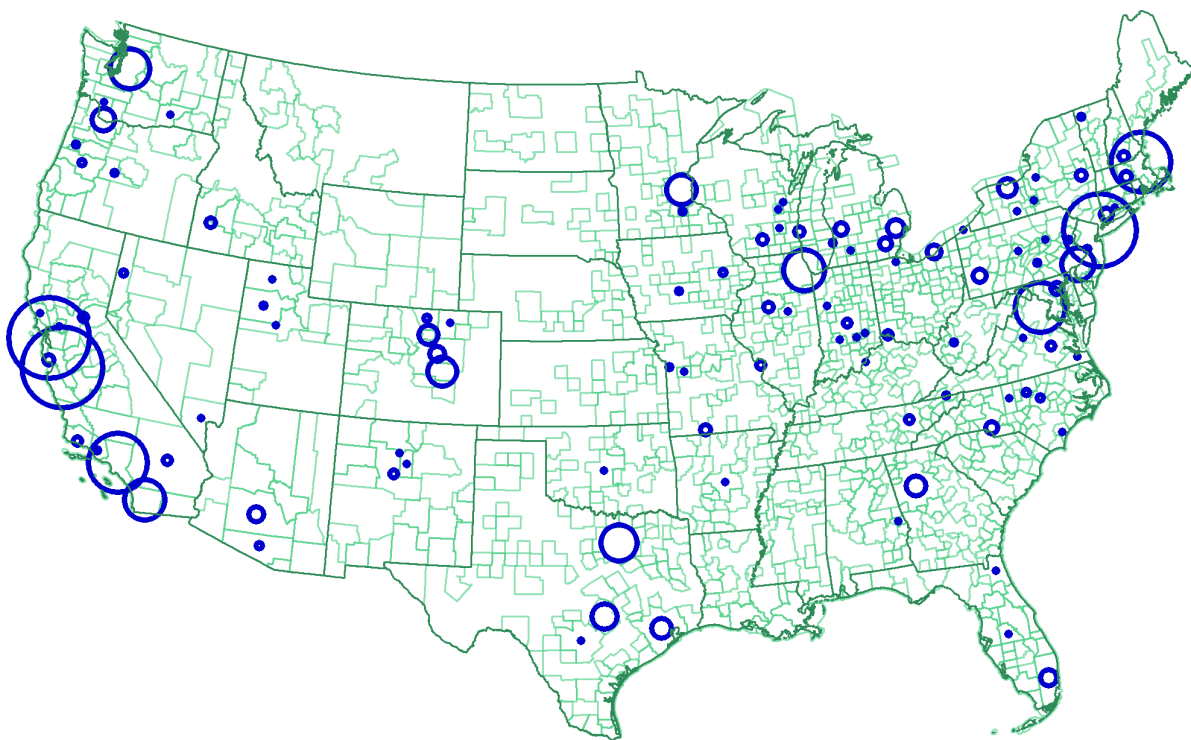
Figure 3-- Share of patents mentioning a new technology filed in Silicon Valley, New York, Seattle and Boston, by technology's economic importance (earnings calls mentions)



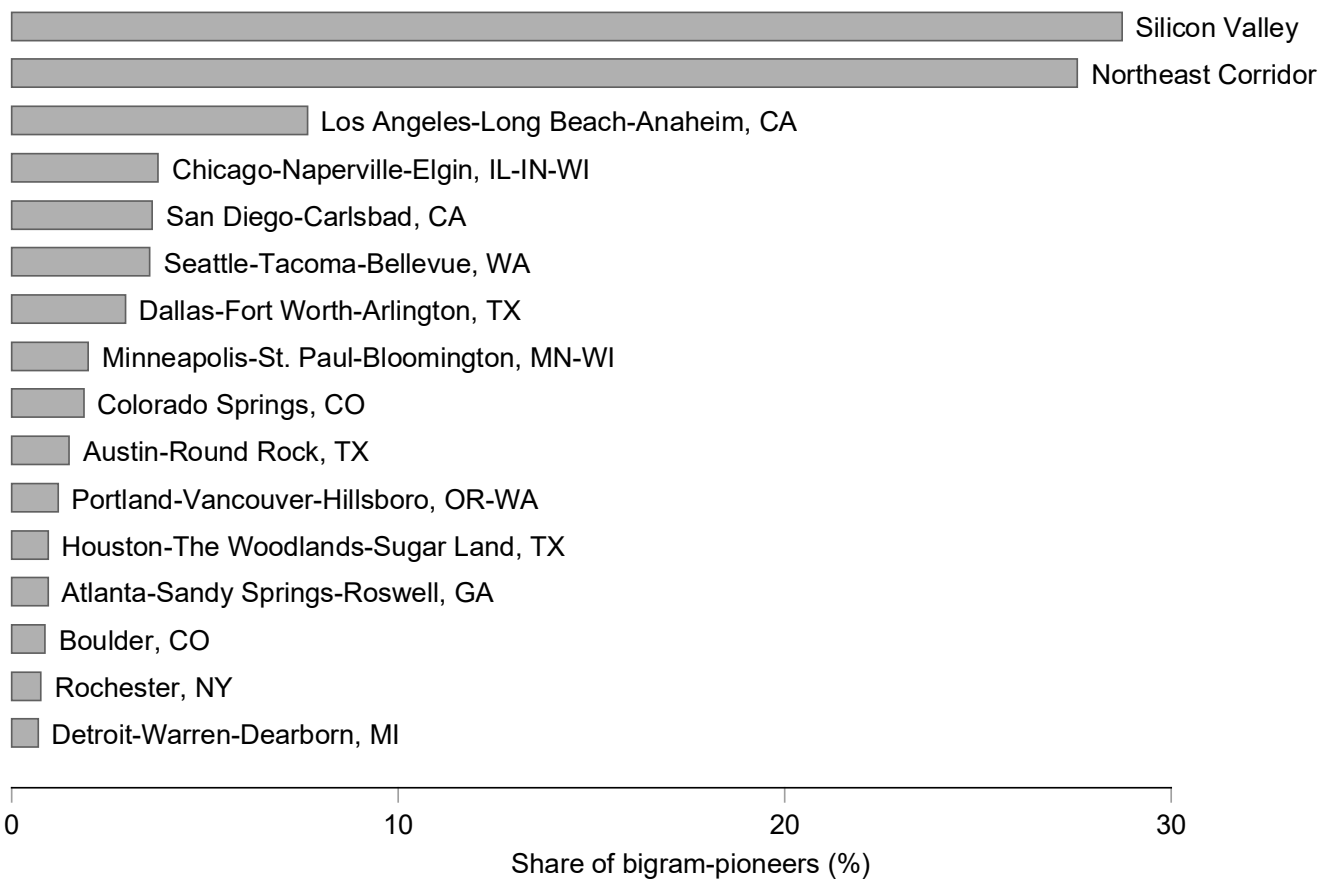
Notes: The figure shows results from a regression of concentration at the bigram-level on indicators for the number of earnings calls that mention a certain bigram, with the respective 95% confidence interval in red. For each bigram, concentration is measured by the share of patenting associated with that bigram in the top five CBSAs (San Jose-Sunnyvale-Santa Clara, CA; San Francisco-Oakland-Hayward, CA; New York-Newark-Jersey City, NY-NJ-PA; Seattle-Tacoma-Bellevue, WA; Boston-Cambridge-Newton, MA-NH). The top five CBSAs are the five regions with the highest number of patents associated with technologies with more than 100 earnings call mentions. Standard errors are clustered by Wikipedia title (technology). Using a F-test, we reject the hypothesis that all coefficients are equal, with a p-value of 0.000 (F-statistic of 42.72).

Figure 4 – Distribution of pioneer locations

Panel A: Pioneer Locations

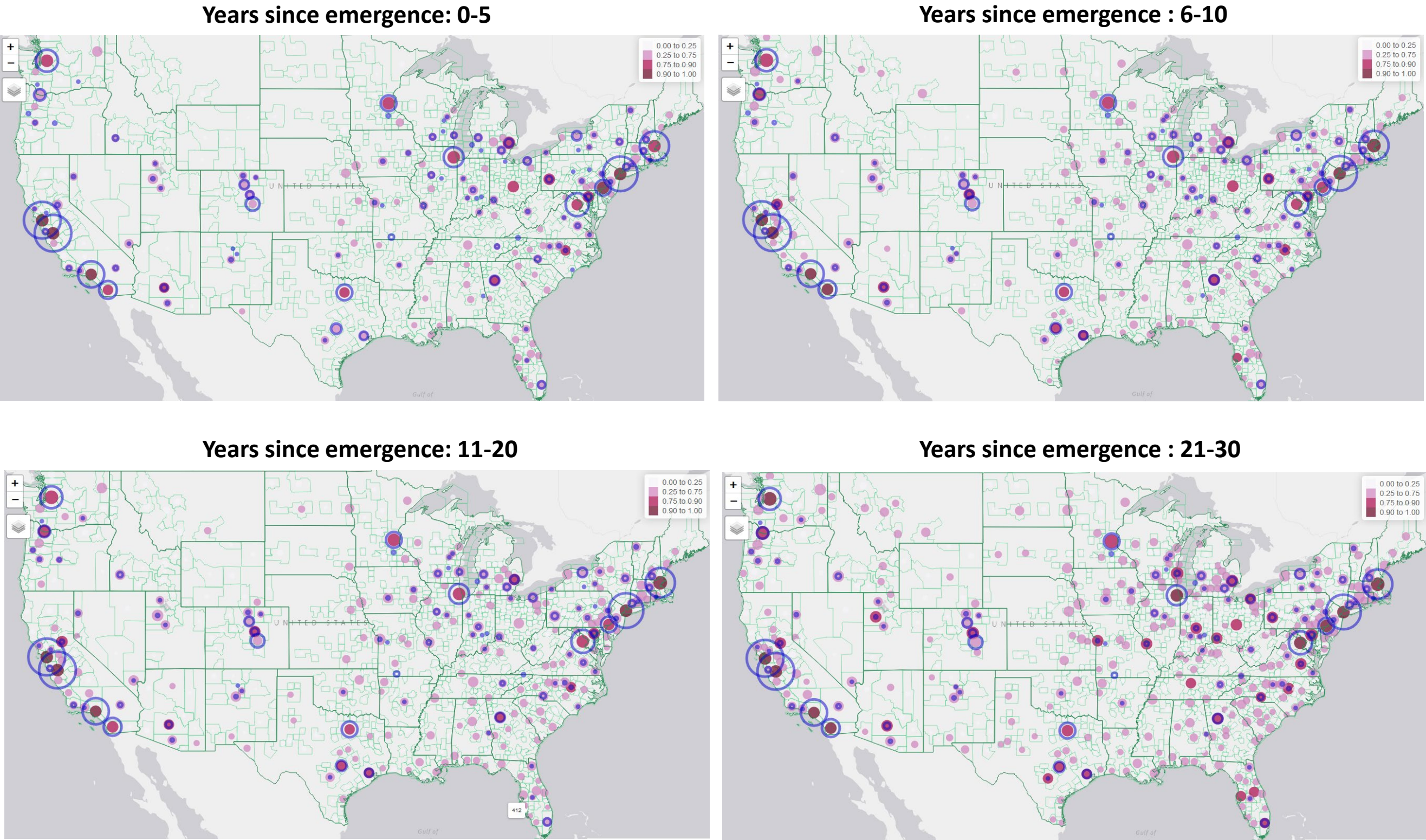


Panel B: Distribution of Pioneer Locations



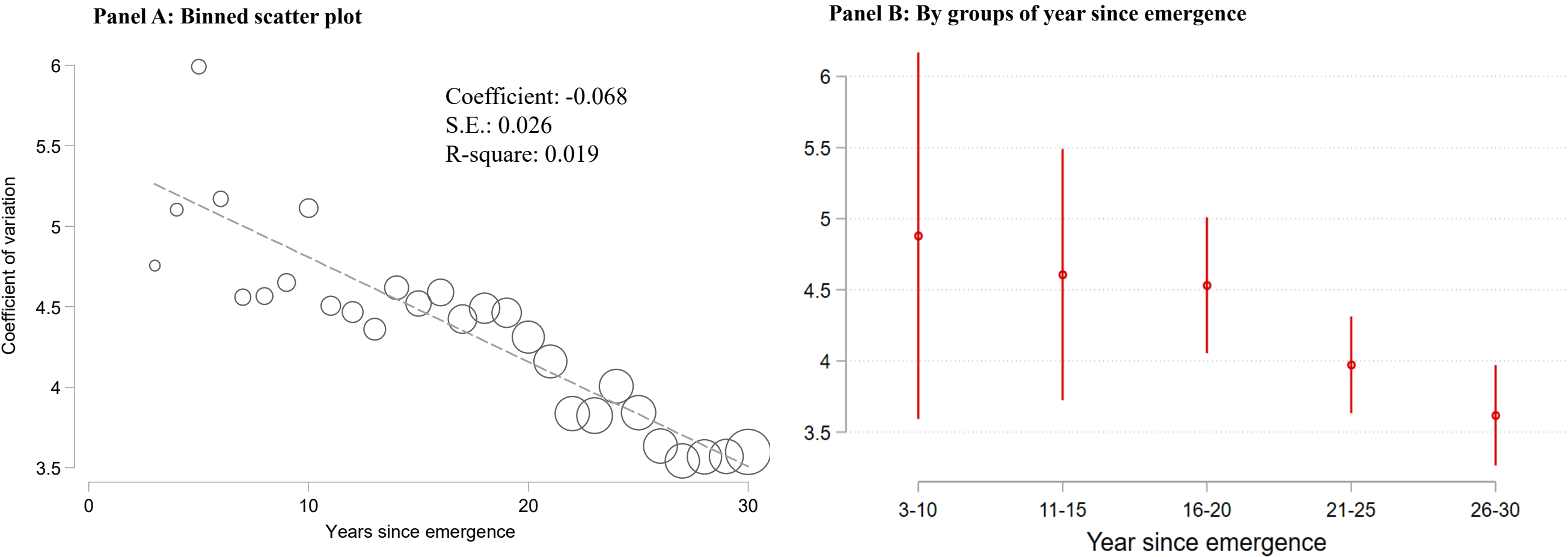
Notes: This figure shows the distribution of pioneer CBSAs. Panel A displays as blue circles CBSAs that are pioneer locations for at least one bigram with more than 100 earnings call mentions. The size of the circles is proportional to the share of technology bigrams for which the CBSA is a pioneer location. Panel B shows a plot of the percentage of technology bigram-pioneer location pairs accounted for by each CBSA, for the top 20 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 5 – Geographic diffusion of technology job postings, by year since emergence



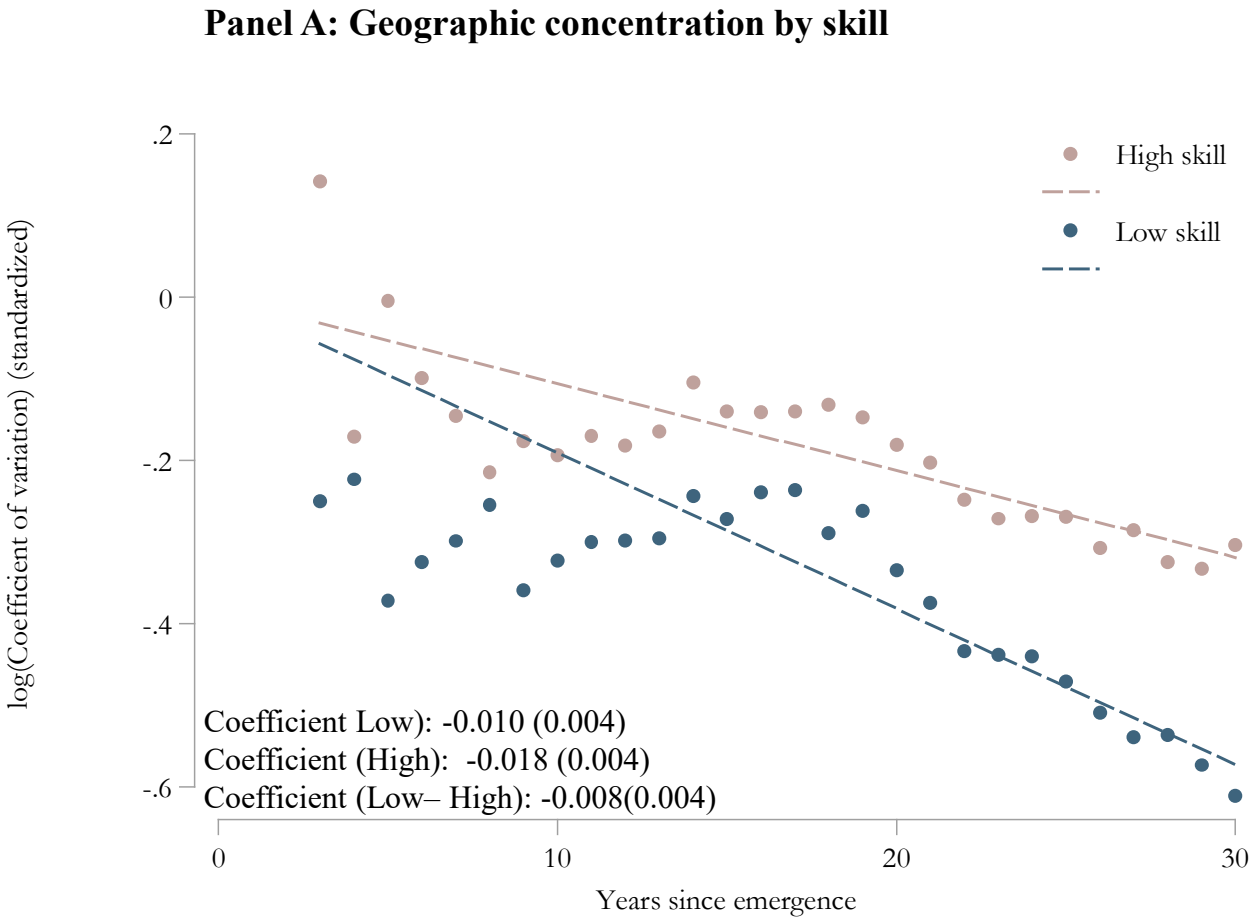
Notes: This figure plots maps with pioneer locations and job postings associated with technology bigrams by year since technology bigram emergence. Pioneer locations are marked with solid blue circles and technology job postings are in solid purple circles. For each CBSA and emergence year, we calculate the share of technology bigrams for which the CBSA records a non-negligible presence of technology jobs ($Normalized\ share_{c,t} \geq 10\%$) 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 6 – Geographic concentration of technology job postings across CBSAs, by year since emergence

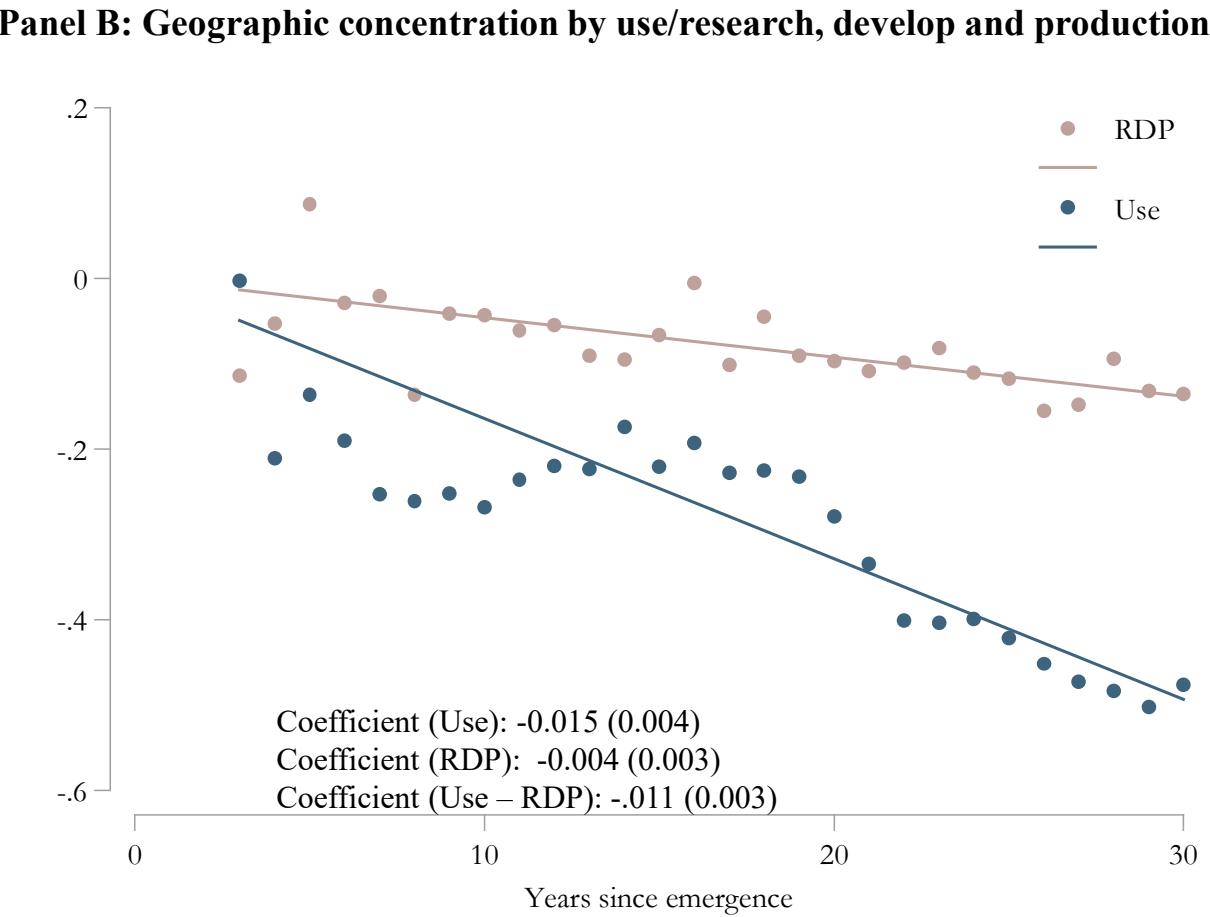


Notes: Panel A shows a binned scatter plot at the technology bigram x year level of the coefficient of variation (CV) of the normalized share of technology job postings over time. We calculate the CV of the normalized share of technology job postings by dividing the standard deviation of $Normalized\ Share_{c,\tau,t}$ across locations c in year t by its mean in year t for each technology bigram τ . 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. Year since emergence is capped at 30 years. Observations in and after the year of emergence are included. Panel B shows the average of the coefficient of variation separately for five age groups (3-10, 11-15, 16-20, 21-25, and 26-30+), also showing standard error bands for these estimates. We only include technology bigrams that appear in at least 100 earnings calls.

Figure 7 – Geographic concentration relative to year since emergence, by skill and type of job posting

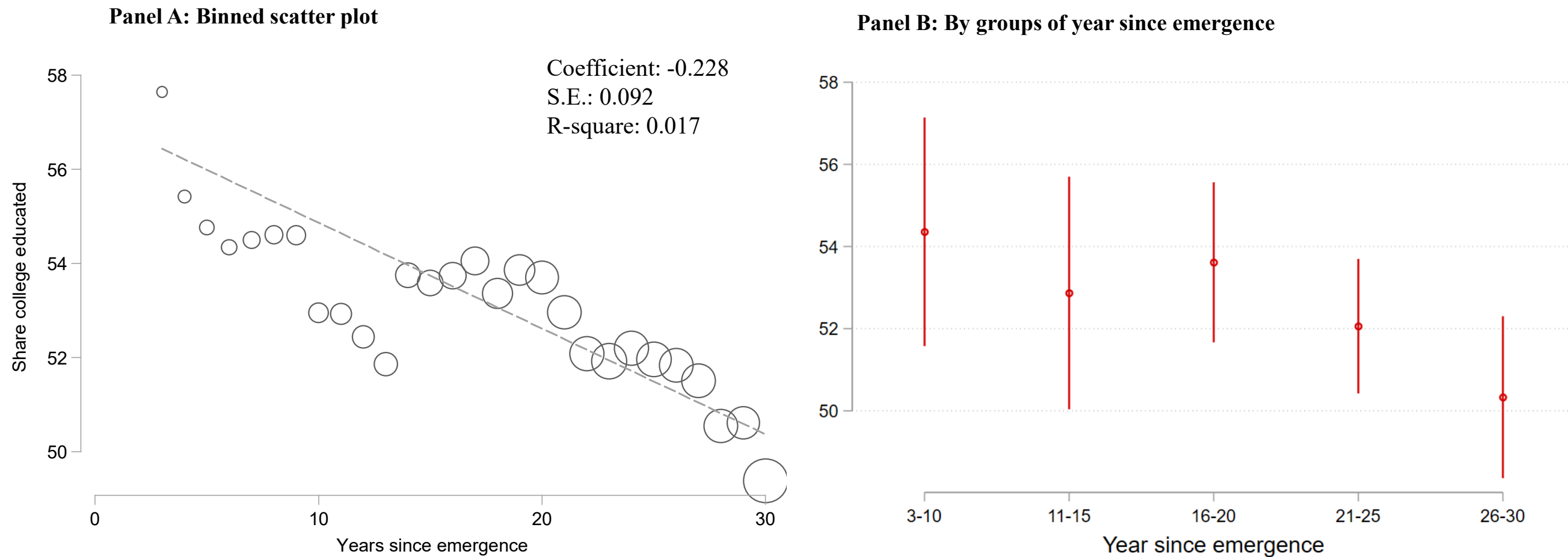


Notes: This figure plots a binned scatter plot at the technology bigram x year x skill category level of the log of the coefficient of variation (CV) of the normalized share of bigram job postings relative to the year since emergence, by skill-level of job posting (high and low). The CV is calculated 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 8 - Share of technology job postings requiring a college education, by year since emergence



Notes: Panel A 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. Panel B estimates the average of the of the same dependent variable separately for five groups (3-10, 11-15, 16-20, 21-25, and 26-30+), also showing standard error bands for these estimates. Year since emergence is capped at 30 years. We only include technology bigrams that appear in at least 100 earnings calls. We only include technology bigrams that appear in at least 100 earnings calls.