

# The Diffusion of Disruptive Technologies

Nicholas Bloom, Tarek A. Hassan, Aakash Kalyani, Josh Lerner, and Ahmed Tahoun<sup>1</sup>

October 25, 2021

**Abstract:** We identify novel technologies using textual analysis of patents, job postings, and earnings calls. Our approach enables us to identify and document the diffusion of 29 disruptive technologies across firms and labor markets in the U.S. Five stylized facts emerge from our data. First, the locations where technologies are developed that later disrupt businesses are geographically highly concentrated, even more so than overall patenting. Second, as the technologies mature and the number of new jobs related to them grows, they gradually spread geographically. While initial hiring is concentrated in high-skilled jobs, over time the mean skill level in new positions associated with the technologies declines, broadening the types of jobs that adopt a given technology. At the same time, the geographic diffusion of low-skilled positions is significantly faster than higher-skilled ones, so that the locations where initial discoveries were made retain their leading positions among high-paying positions for decades. Finally, these pioneer locations are more likely to arise in areas with universities and high skilled labor pools.

**Keywords:** Geography, Employment, Innovation, R&D

**JEL Classification:** O31, O32

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

---

<sup>1</sup> Stanford University; Boston University; Boston University; Harvard University; and London Business School. Bloom, Hassan, and Lerner are affiliated with the National Bureau of Economic Research. We thank audiences at the Applied Machine Learning webinar, Baruch College, Bocconi University, CKGSB, Columbia University, Dartmouth University, Duke, Georgia State, the London School of Economics, Nova Business School, the Toulouse Network on Information Technology, the Royal Bank of Australia, the University of Maryland, the University of British Columbia, the University of Chicago, Yeshiva University, and the NBER Summer Institute and EFG meetings for helpful comments. Special thanks go to Lisa Kahn for sharing data, Bledi Taska for help on BGT data queries, Ben Jones and Chad Syverson for excellent discussions, and to William Hartog, Marcela Carvalho Ferreira de Mello, 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 the Institute for New Economic Thinking, the Ewing Marion Kauffman Foundation, Harvard Business School's Division of Research and Doctoral Programs, the Alfred P. Sloan Foundation, the Toulouse Network on Information Technology, and the Wheeler Institute for Business and Development. Lerner has received compensation from advising institutional investors in venture capital funds, venture capital groups, and governments designing policies relevant to venture capital. All errors are our own.

## 1. Introduction

The development of novel technologies, the degree to which they affect jobs, and the speed with which they spread across regions, firms, and industries are key elements in the study of economic growth, economic inequality, entrepreneurship, and firm dynamics. Many authors have sought to understand whether the benefits from the adoption of new technologies accrue primarily to inventors, early investors, highly skilled users, or to society more widely through, for instance, employment and income growth.<sup>2</sup> Other studies, as discussed below, have explored the geography of the development and diffusion of new technologies.

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 separate those innovations that affect jobs and businesses from those that do not.

In this paper, we make use of the full text of millions of patents and job postings and hundreds of thousands of earnings conference calls over the past two decades to make progress on this challenge. In particular, we develop a flexible methodology that allows us to determine which innovations or sets of innovations (“technologies”) affect businesses, 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 five stylized facts about the development and diffusion of disruptive technologies across space, skill levels, and other dimensions.

The first step of our analysis is to develop a methodology for systematically identifying two-word phrases associated with rapidly diffusing technologies (“technical bigrams”) through a series of systematic rules, whose robustness we verify through various diagnostic tests. To this end, we intersect information from two large corpora of text. First, we use the full text of U.S. patents awarded between 1976 and 2016 to isolate two-word combinations that appear in influential patents but were not commonly used before 1970. That is, we isolate language specific to recent influential innovations. Second, we search for these bigrams in the full text of earnings conference calls held by more than 8,000 listed firms between 2002 and 2020, and identify those technical bigrams that have increased by at least a factor of ten in discussions between firm executives and

---

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

investors during our sample period. This procedure highlights a small set of 305 technical bigrams that describe recent technological advances that have disrupted businesses, in the sense that they have become prominent topics of discussion between the firm’s management and its investors in the last two decades. The top three of these are “mobile devices,” “machine learning,” and “cloud computing.”

To aid interpretation, we then group our technical bigrams into sets of technologies, recognizing the fact that, for example, “cloud computing” and “cloud services” refer to closely-related innovations. This approach partly relies on human judgment, aided by machine learning methodologies. Using this “supervised” process, we identify 29 disruptive technologies, which we use for the main analyses in the paper.<sup>3</sup> Taken together, 28.8% of all citation-weighted patents granted by the U.S. Patent and Trademark Office (USPTO) between 2002 and 2016 are involved with the development of at least one of our 29 technologies. In this sense, our disruptive technologies cover a significant part of recent innovative activity. While we make no claim of completeness, we argue that each of these 29 advances had significant implications for businesses and jobs in the United States in the past two decades.

After establishing our list of new technologies, we then identify patents, earnings calls, and job postings that mention these new technologies. We use patents to identify the locations where each of the technologies was developed and earnings calls to identify exposed firms and the year in which the technology started to feature prominently in the conversations between executives and investors (its commercial breakthrough). We then cross-reference our list of technical bigrams with the full text of online job postings to identify 13 million jobs advertised between 2007 and 2020 that use, produce, or develop our disruptive technologies. These granular data uniquely allow us to track the spread of disruptive technologies along a dimension of crucial importance to policymakers: employment. In particular, we examine the evolution of the number, location, and quality of job postings associated with these new technologies.

The key results of this analysis are as follows.

First, the locations where disruptive technologies are developed are geographically highly concentrated, both within and across technologies. Based on patenting activity ten years prior to

---

<sup>3</sup> In a second, “unsupervised” approach we use all technical bigrams “as is”—i.e., with no human processing. This alternative approach yields qualitatively identical results.

each technology’s commercial breakthrough, we show that the typical disruptive technology in our data emerged from only a handful of urban areas, which housed the majority of early patenting in the technology and the vast majority of its early employment up to the year in which it has its commercial breakthrough. We term these specific urban areas the technology’s “pioneer locations.”

Although 23 of the 50 U.S. states host at least one pioneer location according to this definition, their distribution across technologies is remarkably skewed: a few super-clusters are the birthplace of a surprising number of disruptive technologies in our data. Collectively, locations in California alone host a remarkable 40.2% of our technology-pioneer location pairs. Another super-cluster along the Northeast Corridor from Washington to Boston accounts for additional 24.58%. More broadly, we find that the geographic distribution of patenting related to our 29 disruptive technologies is even more skewed than that of patenting in general.

Second, despite this highly skewed initial distribution, as technologies mature and the number of new jobs related to them grows, they gradually spread geographically. On average, one measure of geographic concentration, the coefficient of variation of the share of local jobs associated with a technology across the 917 core-based statistical areas (CBSAs) in the United States, drops by 24% in the first decade after a technology emerges. We see this pattern of “region broadening” in virtually every technology that we examine.

Third, while initial hiring is concentrated on high-skilled jobs, over time, the mean required skill level of the jobs associated with the technologies declines, reflecting a broadening of the types of jobs that adopt a given technology. For example, the average earnings associated with job postings in a given new technology drops by about 15% within the first decade, falling from \$70,468 per year to \$60,608 per year on average, a drop of about one thousand dollars per year (all figures in 2015 dollars). This pattern of an increasing share of low-skilled jobs that begin to use or produce a given technology holds within most (though not all) of our disruptive technologies.

Fourth, region and skill broadening interact: Low-skill jobs associated with a given technology spread out across space significantly faster than high-skill jobs. Our estimates suggest low-skilled jobs that use or produce new technologies are almost fully dispersed geographically within 20 years. For example, as technologies like the smart phone, cloud computing, and electric cars

mature, the lower-skilled jobs associated with these – salespeople, technicians, repair specialists, etc. – spread across the United States at a fast clip.

Fifth, despite region and skill broadening, disruptive technologies appear to yield long-lasting benefits for their pioneer locations, particularly when it comes to high-skill employment. Our estimates imply it takes almost 40 years for high-skilled job postings to fully disperse from their original pioneer locations. Perhaps not surprisingly, these pioneer locations tend to be located around universities and in areas with more educated populations. Thus, regions with strong local education, research institutions, and universities appear to benefit from successful disruptive innovation for substantial periods of time.

While the focus of our analysis is on documenting the major stylized facts about the spread of technologies, the granularity of our data also allows us to study the employment dynamics associated with disruptive technologies for individual locations and firms. As an example of such a more micro-focused analysis, we document a case study of the geographic footprints of two large Detroit-based car manufacturers, and how they evolve after the emergence of technologies relating to self-driving cars. In this instance, we show that both large incumbents shifted significant numbers of job postings relating to self-driving cars towards the technology’s pioneer locations (particularly Silicon Valley) and away from their traditional hub in Detroit. We speculate that this kind of “re-homing” of established firms may form part of the reason for the long-lasting hiring advantages of pioneer locations.

In the final part of the paper, we look at the generality of our results by studying the diffusion of disruptive technologies across firms, industries, and occupations. We show that similar patterns as discussed before predominate. While technology-related job postings spread out over time, the original firms, industries, and occupations associated with the development and early employment in the technology retain an advantage over time. Generally, we find a faster spread across locations and firms than industries and occupations.

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 those jobs that get filled *ex post* (hiring). Second, by its very nature, our data speaks to job openings relating to novel technologies, but not the possible destruction of existing positions as a result of the diffusion of these technologies. A third concern

is what Merton (1968) termed in his work on citations in scientific papers “obliteration by incorporation”: in our context, when a technology becomes so widely diffused that it is no longer mentioned specifically in earnings calls or job postings.

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, Levy, and Murnane, 2003) to broadband (Akerman, Gaarder, and Mogstad, 2015) to robots (Acemoglu and Restrepo, 2020) to artificial intelligence (Agrawal, Gans, and Goldfarb, 2019; Webb, 2020).<sup>4</sup> A second strand focuses on specific innovations during important historical episodes. Examples include studies of hybrid corn (Griliches, 1957), electrification (Goldin and Katz, 1998), threshing machines (Caprettini and Voth, 2020) and encyclopedias (Squicciarini and Voigtlander, 2015). Mokyr (1992) and Gordon (2016) trace out the impact on economic development and the standard of living of a range of great inventions. Both of these classes of studies use technology- and industry-specific approaches to measure the diffusion and impact of individual technologies. A third strand examines the impact of technological progress more generally on the labor market, including inputs like research and development spending (e.g., Berman, Bound, and Griliches, 1994; Machin and Van Reenen, 1998; and Aghion et al., 2019) and outputs like computerization (e.g., Krueger, 1993; Autor, Katz, and Krueger, 1998; Michaels, Natraj, and Van Reenen, 2014). We contribute to this literature by providing a flexible methodology to systematically isolate those innovations that have a large impact on firms and labor markets, and to track their spread across firms, industries, occupations, and jobs requiring different skill levels. Aside from the 29 disruptive technologies we identify in this paper, variants of our approach could also be used to study the adoption and spread of some of the other specific innovations highlighted by this literature.

A second broad literature examines regional development, in particular questions relating to the mechanisms behind the continued advantage of pioneer locations. A number of papers have highlighted persistent advantages in entrepreneurship (Glaeser, Kerr, and Kerr, 2015), and

---

<sup>4</sup> This work is related to Comin and Hobijn (2004, 2010) and their associated work. Their 2010 paper, for instance, looks at the diffusion of 15 technologies across 166 countries, employing a variety of measures of technological utilization. The rich data that we are able to exploit allow us to analyze (albeit for one nation and a much shorter time period) the interactions between innovation and employment at the firm level on a temporal and regional basis.

innovation (Moretti, 2019) that certain urban areas enjoy, and highlighted mechanisms such as employee mobility across new ventures (Gompers, Lerner, and Scharfstein, 2005) and localized knowledge spillovers (e.g. Jaffe, Trajtenberg, and Henderson, 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 over time, and show there is a general relationship between successful innovation, early employment in a given new technology, and the long-term advantage that these locations preserve in high-skill employment. We hope that future work will use our granular data to study in detail the anatomy and evolution of technology hubs.

Third, our work relates to a broader literature on the diffusion of new technologies. Since the pioneering work of Griliches (1957), the diffusion process has long been understood by economists to be a gradual one. While broader sociological and organizational literature has examined the barriers to innovation, recent work in economics has focused on understanding the importance of supply and demand factors on the speed of diffusion (e.g., Popp, 2002; Acemoglu and Linn, 2004; Greenstone, Hornbeck and Moretti, 2010; Moser, Voena, and Waldinger, 2014; Moscona, 2020; Arora, Belenzon, and Sheer, 2021). Despite this interesting work, Hall’s (2004) characterization of the study of diffusion as “a somewhat neglected one in the economics of innovation” still remains a fair observation. Our contribution is to provide an assessment of the rate at which disruptive technologies spread across locations, firms, occupations, and industries.

Finally, our work adds to a growing literature in economics and related fields using text as data. In particular, a number of recent papers have used newspapers, patents, and firm-level texts to measure concepts that are otherwise hard to quantify using conventional data sources. (Examples include Baker et al., 2016; Hassan et al., 2019, 2021; Handley and Li, 2020; Sautner et al., 2020; Bybee et al., 2019; and Kelly et al, 2021). We focus primarily on the full text of job postings, which has received relatively less attention.<sup>5</sup> Our work adds to this literature by introducing a methodology to jointly analyze the full text of patents, earnings calls, and job postings.

The remainder of this paper is structured as follows. In Section 2, we describe the construction of our data. In Section 3, we present our region-broadening and skill-broadening results. In addition, we examine the differential patterns across geographic regions. In Section 4, we examine the

---

<sup>5</sup> A notable exception is the work by Veldkamp and Abis (2021) who use job descriptions to identify financial analysis positions that leverage machine learning.

diffusion of disruptive technologies across three other dimensions: industries, occupations, and firms; and investigate a potential mechanism for region broadening: firm rehoming towards pioneer locations. Section 5 shows a number of additional robustness checks.

## 2. Data Construction

In this section, we describe our text-based approach to identifying and tracking the spread of disruptive technologies. We first intersect the full text of patents with that of earnings conference calls to identify those keywords describing new technologies that have increasingly appeared in conversations between investors and executives at listed firms. We then audit and group these keywords into 29 distinct disruptive technologies and track their use across patents, earnings calls, and job postings. Our objective is to (a) build a firm-quarter level measure of technology exposure, (b) use this measure to pinpoint when a given technology starts affecting businesses; (c) identify which locations and firms pioneered early patents in a given disruptive technology; (d) create a measure of technology adoption at the job-posting level, and (e) aggregate this measure to track the spread of disruptive technologies across jobs posted in different regions, occupations, firms, and industries.

### 2.1. Phrases unique to novel and influential patents

As a first step, we want to identify influential technologies in as systematic a manner as possible. We begin by examining U.S. patent filings. Patents are an attractive starting point for our analysis for two reasons. First, they are by definition novel, particularly when we focus on the most influential patents. Second, they must describe their technology and (at least some) key ways in which it is applied.<sup>6</sup> We focus solely on patent awards by the USPTO: because of the importance of the U.S. market, inventors worldwide typically file important discoveries with the USPTO.<sup>7</sup>

In order to obtain set of bigrams associated with novel technologies, we collect all utility patents awarded by the USPTO to either U.S. assignees or inventors between 1976 and 2016, a total of approximately three million awards. From the text of these patents (abstract, summary, claims, and

---

<sup>6</sup> This requirement is stipulated in the legal concept of “reduction to practice,” 35 U.S.C. 112(a).

<sup>7</sup> About half of all patent applications to the USPTO are filed by residents of foreign countries (USPTO, 2020). This pattern reflects the fact that patent protection in a given nation depends critically on having a patent issued in that specific nation. Important discoveries (the focus of our analysis) are therefore disproportionately likely to be filed in major patent offices world-wide (Lanjouw, Pakes, and Putnam, 1998).



background description), we remove stop words (such as “of,” “the,” and “from”) and decompose the remaining text into about 17 million unique two-word combinations (“bigrams”). 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, Wang, and Lee, 2002; Bekkerman and Allan, 2004).

To reduce the 17 million bigrams appearing in patents to a more manageable number, we next take steps to isolate those bigrams that are associated with *novel* and *influential* inventions. First, we focus our attention on bigrams associated exclusively with *novel* inventions by dropping “non-technical” bigrams that were in common use long before the emergence of our disruptive technologies. To this end, we select all text dating prior to 1970 from the Corpus of Historical American English (COHA), a representative sample of text constructed by linguists from prominent fiction and non-fiction sources (Davies, 2009) that reflects everyday use of English up to 1970. We then remove any bigram appearing in this source (for instance, “of the” or “equipment used”) from our list of bigrams obtained from patents, leaving us with 1.5 million exclusively “technical” bigrams.

Second, to identify bigrams associated with *influential* inventions in the remaining list, we collect patent citations for all the patents that mention these bigrams between 1976 and 2016 and normalize the citations to each patent by the mean within each technology class and application year.<sup>8</sup> We then retain only those technical bigrams that cumulatively obtain at least 1000 normalized citations. After these eliminations, we have a list of 35,063 “technical” bigrams associated with influential patents between 1976 and 2016.

Next, we focus on which of our technical bigrams figured increasingly into the business discussions of firms, to gauge the extent to which each innovation changed or disrupted how firms operated. Here we use earnings conference calls from publicly listed firms.

## 2.2. Use of technical bigrams in earnings calls

Quarterly earnings call transcripts consist of two sections: a presentation by management (typically the chief executive and/or financial officer(s)) and then questions posed by investors and analysts

---

<sup>8</sup>Citation rates vary considerably over time and across technology classes. Lerner and Seru (forthcoming) document this heterogeneity and the biases that can result from failing to correct properly for these differences.

with answers provided by the executives. These calls have been shown to be indicators of some of the most important issues facing these organizations (Bushee, Matsumoto, and Miller, 2003; Matsumoto, Prank, and Roelofsen, 2011; Hassan et al., 2019, 2021).

We tabulate the bigrams in 321,373 conference calls held by 11,905 publicly held companies and compiled by Refinitiv EIKON (formerly Thomson Reuters) between 2002 and 2019. Through this examination, we eliminate about 43% of our technical bigrams from the patents that are never mentioned in these calls.

We then trim the remaining bigrams in two ways. First, we require that they appear in more than 100 transcripts, to focus only on economically important bigrams associated with innovations that became major topics in earnings discussions. Second, we require these are increasing in their incidence in earnings calls over time, to focus on technologies that disrupt businesses, in the sense that they become a growing topic of conversation between executives and investors during our sample period. We thus define a “disruptive” technology simply as one that takes up an increasing amount of airtime in the earnings conference calls of listed firms. In our baseline specification, we keep technical bigrams which appear at least ten times as frequently in their peak year as in the first year of the earnings call data in 2002.<sup>9</sup> After these steps, we end up with a short list of only 305 technical bigrams describing technologies which are widely used and rising in importance, which we label as disruptive technologies.

Table 1 shows the 30 technical bigrams most frequently appearing in earnings calls. It shows that our simple two-step approach of cross-referencing bigrams from influential patents with those featuring increasingly in business discussions clearly identifies some of the major disruptive technological advances of the past two decades. The first four bigrams on the list are “mobile devices,” “machine learning,” “cloud computing,” and “cloud services.” Other top-ranking bigrams on the list include “social networking” and “smart grid.”

In order to obtain a coherent set of technologies from our 305 bigrams, we take two approaches. In the first, described in detail below, we manually group the 305 bigrams into a set of technologies, recognizing the fact that, for example, “cloud computing” and “cloud services” refer to closely related innovations. We apply a number of further refinements, allowing us to quantify

---

<sup>9</sup> Bigrams that do not appear at all in any call held during the 2002 calendar year automatically meet this criterion.

the spread of specific technologies along a variety of dimensions. This approach inevitably relies on human judgment, aided by machine learning methodologies. This “supervised” approach is the basis for the analyses presented in the main body of the paper. We describe it in detail below.

An alternative “unsupervised” approach is to use all 305 bigrams “as is”—i.e., with no human processing. We show in Section 5 that all of our main results are robust to this approach, both qualitatively and quantitatively. In this sense, all human intervention from this point on serves the purpose of measuring the spread of specific technologies and making our results more easily interpretable, but has no bearing on the validity of our main stylized facts about disruptive technologies as a whole.

Before describing our supervised approach in detail, it is worth commenting on our perhaps remarkable finding that of the 35,063 “technical” bigrams that are unique to novel and influential inventions only a few hundred go on to “disrupt” a large number of conversations between executives and investors. First, the distribution of technical bigrams across earnings calls is highly skewed (the median bigram is mentioned in only one call), so that only 2,181 technical bigrams feature in more than 100 earnings calls – our threshold for our notion of economically important innovations. Of these, only our 305 technical bigrams increase in frequency of use by factor ten during our sample period – our threshold for considering it “disruptive.” Interestingly, most of these bigrams (235) also increase by a factor of 100, so that varying this threshold has only a very modest effect on the bigrams selected. For example, under this more stringent notion of “disruptive” technologies, the three most frequent technical bigrams are “machine learning,” “cloud computing,” and “cloud services.” Appendix Figure 1 and Appendix Tables 1 and 2 show more systematically how the top and bottom technical bigrams change with other reasonable choices for both cutoffs (requiring a minimum of 80 to 120 mentions or an increase in frequency by factor 5 to 100).

Generally, as both cutoffs become more stringent, our procedure isolates technologies that are more broadly impactful and disruptive. As we loosen the cutoffs, we pick up a moderate number of additional technologies (such as laser welding), but also additional noise. More importantly, as we show below, all of our main stylized facts remain unchanged with these alternative sets of “unsupervised” keywords.

### 2.3. A “supervised” approach to defining specific technologies

In our “supervised” approach, we take four additional steps to ensure we accurately track the spread of specific technologies. First, we eliminate those bigrams (from the list of 305 bigrams) that, in our reading, do not clearly and unambiguously reflect specific technological advances. This approach allows us to eliminate bigrams that refer to problems that have recently become more salient for both inventors and executives, but are not technological solutions, such as “carbon footprint” or “power outage.” Similarly, we drop bigrams referring to older technologies, such as “smart grid,” which refers to a technology that has been available since the 1980s but is enjoying renewed interest in recent years, and “nand flash” (flash memory), which had a surge of references when a global supply issue occurred. We also drop any bigram that is vague or refers to multiple innovations, such as “flow profile,” which may refer interchangeably to a genomic flow or firms’ cash flows, and “digital channel,” which can refer to interchangeably to digital marketing or digital transmission. At the end of these eliminations, we retain 105 bigrams that, in our reading, clearly and unambiguously reflect specific disruptive technological advances.

We then manually pair each of the remaining bigrams with a definition sourced from Wikipedia and form 29 groups of bigrams (“technologies”) that each refer to a specific technological advance defined in this source. For example, the bigrams “mobile devices,” “smart phones,” and “mobile platform” all refer to “smart devices,” which Wikipedia defines as “an electronic device, generally connected to other devices or networks via different wireless protocols.” Appendix Table 3 lists the definition used for each of our 29 technologies.<sup>10</sup>

A general concern with our approach of intersecting different corpora of text to measure the spread of technologies is that the phrases used by executives to characterize new technologies may never appear in patent awards, leading us to under-count mentions of some technologies relative to others. To explore this possibility, we use an embedding vector algorithm (word2vec, developed by Mikolov et al., 2013 and used previously by Hansen et al., 2021 and Atalay et al., 2020), which we trained specifically on our set of earnings calls. This algorithm trains a neural network to map each bigram as an embedded vector into a multidimensional space, such that bigrams which are

---

<sup>10</sup> We explored automating this grouping procedure. For instance, we experimented with clustering two bigrams into a group if the average similarity from the patent and EC embedding vectors was more than 70%. This gave a similar grouping as when using our human judgement. However, when differences arose between the automated and human approaches, we generally preferred the results based on our human judgement, so we used the latter as our preferred approach. For example, in the automated approach, “virtual reality” and “augmented reality” were clustered together with “machine learning” and “neural network,” while in our human approach we split these into two technologies: “virtual reality” for the first two and “machine learning/AI” for the second two.

used in common contexts are located closer in this space. We use the trained algorithm to derive a set of bigrams closest in terms of cosine distance to a given group of technical bigrams. For each of our 29 groupings, the algorithm suggests a list of “proximate” (similar) bigrams. For example, the most proximate bigrams to those in the technology grouping “artificial intelligence” are “machine learning” and “deep learning.” From this list, we then add to the bigrams forming each technology those that, in our reading, also clearly and unambiguously describe the technology in question.

At the end of the process, we also wish to ensure that our audited list of bigrams correctly identifies postings of jobs involved with using or producing a given technology. To this end, we performed an iterative human audit where a team member went through 3,460 randomly sampled excerpts of the text from job postings (covering at least 10 unique postings per bigram). He or she classified the snippet into true positive and false positive categories, along with suggestions regarding new keywords discovered and how the accuracy of the existing keywords could be improved.<sup>11</sup> We retained only bigrams that, according to our reading, unambiguously reflected discussion of the technology in question at least 80% of the time. For example, we find that the bigram “automated car” rarely refers to the “Autonomous Car” technology but instead to automated car washes. Appendix Table 4 shows this human audit process in detail for an example.

Following these additions and subtractions, we obtain a list of 221 audited technical bigrams associated with our 29 disruptive technologies.

Table 2 lists the 29 disruptive technologies from our supervised approach and the associated number of Burning Glass job postings in which associated bigrams appear (see the discussion in the next section). In addition to the major innovations already mentioned above, they include well-known green technologies (“Solar Power,” “Hybrid Vehicle/Electric Car”) and process innovations, such as “3D Printing,” “Fracking,” and “Machine Learning,” but also less well-known technical and medical advances (e.g., “Millimeter Wave,” a novel band of radio frequency, and “Antibody Drug Conjugates,” a class of drugs used for the treatment of cancer).

---

<sup>11</sup> As an example of a false positive, an ad for a truck driver asked “do you hold a current Class A or B commercial driver’s license with an air brake endorsement? ... do you enjoy playing video games or computer games with a joy stick? are you good at backing up in tight spaces?” The second question led the job to be (incorrectly) classified under “electronic gaming.”

Taken together, they cover a broad range of new methods and consumer applications. While we make no claim of completeness -- other methods might well yield different groupings and definitions of technologies – we show below that each of these 29 advancements had significant implications for businesses and jobs in the United States.

Table 3 lists all bigrams used for the first ten of our technologies in alphabetical order. Appendix Table 5 provides the full set of bigrams used for all technologies.

## 2.4. Burning Glass Job Postings

Burning Glass (BG) aggregates online job postings using “spider bots” from online job boards (such as indeed.com), employer websites (such as stanford.edu), and other sources into a machine readable, de-duplicated database. From Burning Glass, we employ two datasets. The first is a standardized dataset (used recently by Hershbein and Kahn, 2018; Demming, 2020; and Atalay et al., 2020) where each de-duplicated job posting is geo-coded and assigned to a Standard Occupational Classification (SOC) code, a United States government system of classifying occupations, and a North American Industry Classification (NAICS) code.<sup>12</sup> The second dataset has thus far received less attention by researchers. It contains the raw unprocessed text of the job postings, which we use to assign exposure to our technologies.

We use these data from BG for all available years, 2007 and 2010-2020, a total of roughly 200 million job postings. We show below that all of our main results are robust to dropping the 2007 vintage from the sample.<sup>13</sup>

We associate each posting with a skill level, location, industry, and firm as follows (for details, see Appendix 1.4):

- Skill level: We construct a skill level for each six-digit SOC code (the most detailed level) from BG by measuring the share of persons with a college degree, the share of persons with a PhD, the average wage, and the average years of schooling in the American

---

<sup>12</sup> We make extensive use of the former, which are available for 80% of all postings. Industry classifications are available for a more limited 41% of postings. We use these only in our calculations in Section 4. The strings with firm names are available for 66% of all postings.

<sup>13</sup> BG’s efforts to compile job postings data were interrupted by the 2008-09 recession.

Communities Survey (ACS 2015 release), using respondents reporting their occupation as in that six-digit SOC code.<sup>14</sup>

- Location: We use the county names provided by BG to assign job postings to a core-based statistical area (CBSA), a U.S. government-defined geographic area that consists of one or more counties (or equivalents) with an urban center of at least 10,000 people, plus adjacent counties that are socioeconomically tied to the urban center. In total, the dataset includes 917 CBSAs.
- Industry: We allocate a job posting to an industry using the four-digit NAICS code provided by BG.<sup>15</sup>
- Firm: BG reports an employer string for about 60% of their job postings. In order to match these employer strings to firms, we extend the methodology of Autor et. al. (2020) as follows: We search for the employer string (lower case and only letters a-z) on Bing.com, and collect the top five search results. We identify pairs of employer strings as the same firm if they share at least two out of top five search results. We then cluster together all employer strings that have at least two results for the same firm, and associate them with that firm.

## 2.5. Constructing the Exposure Measures

Using these data, we then construct measures of exposure to the set of technologies for job postings, earnings calls, and patents using the following rule:

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

where  $D_{i,\tau}$  is the set of bigrams contained in a job posting/earnings call/patent that was posted/held/filed at time  $t$  and  $b_\tau$  is a bigram associated with a technology  $\tau$ . A document is thus classified as exposed to a technology if it contains a bigram associated with the technology.

Though we use the same terminology to refer to exposed job postings, earnings calls, and patents, it is worth emphasizing that these three types of exposures naturally have different interpretations.

---

<sup>14</sup> For SOC codes in job postings where we do not find any persons surveyed in the ACS, we match them to the closest available SOC code in the ACS. For example, data for SOC Code 38-1967 was not available, so we match it to 38-1960. In total, the dataset includes 837 SOC codes.

<sup>15</sup> NAICS codes typically have six nested levels; the four-digit level is referred to as “industry group.”

Patents that mention one of our technologies are, of course, in some way related to the development of the technology. Appendix Figure 2 provides an example of a patent concerned with object recognition, which mentions the bigram “object recognition” (a keyword associated with our “Computer Vision” technology) 52 times. Similarly, firms exposed to a given technology might be involved in producing or using a given technology, but they may also compete with or be disrupted by the technology. Appendix Table 6 gives text-based examples of these different kinds of firm-level exposures measured from earnings calls.

Most importantly, the vast majority of job postings that mention a given technology advertise jobs that either produce or use a given technology. Figure 1 and Figure 2 provide examples of two illustrative Burning Glass job postings exposed to AI and solar technology, respectively. The first is for an applied research scientist and requires “knowledge of *machine learning*, *neural networks*, and *deep learning*” – all bigrams we associate with the “Artificial Intelligence” technology. The second is for a solar panel installer, and lists as part of the job’s responsibilities “install the racking system and *solar panels*.” Further down, this posting also contains another, more problematic mention of the same technology in the context of the company description, not the job itself.

To investigate the context of technology exposure in job postings more systematically, one of our team members went through 100 randomly sampled job postings for each of our 29 technologies. He or she classified them into two sets of categories, whether technology exposure in the posting referred to 1) either the overall company description or the specific task of the job in the posting, and 2) either the use or the production of the technology.

Appendix Table 7 summarizes the findings from this analysis. In Panel A, we report that in 80% of the postings, the technology mentions refer specifically to the job task (as in Figures 1 and 2). These are split about half and half into the use and the production of the technology. An example of produce would be “*You will be designing the graphics module for our **virtual reality** training system*” while an example of use would be “*The role will involve assisting customers and selling tickets from your **smart tablet** in the entrance of the cinema*”. During the audit, we also noted that company descriptions are usually in the beginning or towards the end of job postings. For this reason, we disregard any technology mentions in the top and bottom 50 words of each job posting. This procedure increased the rate of capturing specific job-related tasks associated with the technology to 91% in our human audit. An additional 4% of mentions were unspecific (for



example, mentions of these technologies being available in the workspace), and only 5% referred to the company but not the job.

In total, we find our 221 technical bigrams mentioned in 13 million job postings, where on average each bigram appears in 59,013 postings. To put this number into perspective, it is useful to compare this frequency with the frequency of other “non-technical” bigrams often used by investors and executives in earnings conference calls. As documented in Appendix Tables 8 and 9, we here reverse our methodology and, instead of selecting bigrams that appear in both patents and earnings calls, select those that appear in earnings calls but not in patents. We find that the top 221 most frequent non-technical bigrams from earnings calls are on average mentioned in only 142 job postings. That is, our technical bigrams are mentioned four hundred times more frequently in job postings than other language frequently used by investors and executives, already suggesting that our 29 disruptive technologies indeed had a large impact on the U.S. labor market.

Having constructed our document-level exposure measures, we next aggregate over various documents  $D$  (job postings, earnings calls, and patents) to construct measures at the region, sector, occupation, and firm levels:

$$share\ exposed_{a,\tau,t} = \frac{\sum_{i \in a,t} 1\{b_\tau \in D_{i,t}\}}{\sum_{i \in a,t} 1\{D_{i,t}\}} \quad (2)$$

where  $a$  may be a region sector, region, occupation, or firm and  $t$  is time. To illustrate, Appendix Table 10 shows a list of top occupations exposed to one of our technologies, virtual reality. Appendix Tables 11, 12, and 13 provide a shorter list of the most exposed regions, occupations, and industries for each technology.

### 3. Diffusion across Regions and Skill-levels

We first seek to understand the overall patterns in the diffusion of these 29 technologies. The analysis suggests that job postings referring to given technologies grow in tandem with references in earning calls; and that over time, hiring moves from a sharp focus on high-skilled jobs to a much broader intake of workers with lower skills.

Figure 3 takes a first look at the diffusion of disruptive technologies. The 29 images plot measures of activity in job postings and in earnings calls on an annual basis for each technology. The red line denotes the percentage of firms in earnings calls that mention the given technology. In some

cases, such as touchscreen and RFID, the number of mentions climb and then fade, presumably reflecting the increasing ubiquity, and hence the declining competitive relevance, of the technologies for firms. In others, such as 3-D printing and artificial intelligence, there is a steady climb over time.

In each plot, we mark the year in which the technology became economically significant, which we henceforth refer to as the “emergence year.” To compute this, for each of our technologies, we calculate the maximum of the “percentage of earnings calls” time series graphed in Figure 3. We define the emergence date to be the year in which the time series first attains at least 10% of this value. Appendix Table 14 lists the emergence date for each technology, along with an alternative definition using the time series of the share of patents exposed to the technology. All of our main results are unchanged when we use this alternative definition.<sup>16</sup> Appendix Table 3 lists each technology, its definition as discussed above, and a suggested contemporaneous event around the year of emergence of the technology.

The second series in Figure 3, denoted with gray dots, indicates the share of positions in Burning Glass that mention a given technology (the size of the dots scale with the number of jobs posted). While in some cases a given technology continues to be important in hiring even after its mentions in earning calls drop off (e.g., GPS technology), in general, the two series are quite closely correlated. The correlation coefficient between them across the figures is 0.81. The close tie between these series helps validate the reasonableness of our empirical methodology: when a technology becomes more commercially relevant for firms, it also becomes more relevant for jobs.<sup>17</sup>

Consistent with this pattern, we also find that more extensive discussions of a technology in earnings calls correlate strongly with more patenting activity in that technology. Appendix Figure 3 shows the share of firms exposed to each technology (in red-solid), and the share of citation-weighted patents (normalized by the average number of citations within each technology class and

---

<sup>16</sup> This alternative definition instead uses the year in which the cumulative number of USPTO patents in that technology attains 50% of its sample maximum. See Appendix Table 19.

<sup>17</sup> Cumulating the number of jobs posted in successive years (in combination with appropriate assumptions on matching and separation rates) shows an increasing stock of technology adoption over time -- a finding reminiscent of Griliches' (1957) S-curves. Appendix Figure 4 shows this pattern graphically for those technologies with an emergence year post 2007.

year) associated with each of our 29 technologies (in black-dashes). Again, the series are highly correlated: the correlation coefficient is 0.80.

### *Region Broadening*

Figure 3 already shows that there is an increase the number of job postings that mention disruptive technologies over time. Figure 4 highlights a related feature: this increasing use in job announcements over time is associated with greater geographic diffusion. To show this, we compute the coefficient of variation in the years after the emergence of a technology (defined as above) measured across locations. More specifically, we create the normalized share of job postings in technology  $\tau$  and year  $t$  for each CBSA-technology-year triple by calculating:

$$Normalized\ share_{cbsa,\tau,t} = \frac{share\ exposed_{cbsa,\tau,t}}{share\ exposed_{\tau,t}}, \quad (3)$$

where the numerator is defined as in (2) and the denominator is the average share of jobs exposed to technology  $\tau$  across CBSAs.  $Normalized\ share_{cbsa,\tau,t}$  thus measures the regional over or underrepresentation of job postings associated with each technology relative to the overall distribution. This normalization allows us to control for the facts that, for instance, Los Angeles, the largest CBSA, will have a large share of job postings of nearly every type and that different technologies may be implemented at very different scales at a given point in time. Appendix Table 15 summarizes the data used in the analysis.

Figure 4 depicts, for each technology and year since emergence, the ratio of the standard deviation and the mean of this measure across CBSAs, also known as the coefficient of variation. The analysis reveals an intriguing pattern: 28 of 29 technologies exhibit a decline in the coefficient of variation over time (the only exception being job postings associated with the “Search Engine” technology). Put another way, although job postings in a given technology are highly regionally concentrated in the early years after their emergence, the geographic distribution of adoption over time becomes more homogeneous.

Figure 4 is corroborated by Table 4, which examines these patterns using a regression framework. Column 1 presents the results of a regression of the coefficient of variation on the years since emergence for an annual panel of technologies, with technology and year fixed effects. Observations are weighted by the square root of the number of job postings associated with a given

technology in each year, in order to give more weight to coefficients of variation that are measured more accurately.<sup>18</sup>

Our preferred estimate in Column 1 shows that the coefficient of variation declines by 0.105 (s.e.= 0.027) per year. The mean coefficient of variation across technologies and years is 4.74. Thus, this estimate implies that the regional concentration of technology job postings declines by 22.1% of the sample average in the ten years after the emergence of the technology.

The remaining columns show the same pattern, using alternative measures of concentration. Column 2 uses the ratio of the normalized share of technology jobs of the top five CBSAs relative to all CBSAs. Column 3 uses the share of CBSAs with a (negligible) normalized share of technology employment of less than 1%. Both variations show concentration significantly decreasing over time.

### *Pioneer Locations*

We next examine the hiring advantage of pioneer locations that excel in initial technology-related inventions. More specifically, we define pioneer locations as those which collectively accounted for 50% of the cite-weighted patent grants associated with a given technology in the ten years before its emergence year.<sup>19</sup> For example, the CBSAs surrounding Trenton (NJ, 21.7%), New York (NY, 11.5%), Rochester (NY, 9.9%), and Los Angeles (CA, 9.3%) are pioneer locations for OLED Display technology because they together accounted for 52.2% of total OLED Display patenting in the U.S. Appendix Table 16 shows the top pioneer location for each of our 29 technologies.

Panel A of Figure 5 shows the geographical distribution of pioneer locations across the United States, where the size of the blue circles is proportional to the share of the 189 technology-pioneer location pairs situated in a given CBSA. Although 23 of the 50 states host at least one pioneer location, the map shows remarkable concentration in this kind of successful innovative activity. Silicon Valley (the San Jose Jose-Sunnyvale-Santa Clara CBSA) and San Francisco were each involved in the development of 23 of our disruptive technologies, followed by New York (21), Boston (18), and Los Angeles (17). Figure 6 shows that collectively, locations in California alone

---

<sup>18</sup> This weighting scheme is for accuracy of our estimates and has no impact on the qualitative results. See Appendix Table 20 for details.

<sup>19</sup> An alternative approach is to define pioneer locations using the regional distribution of a given technology's job postings prior to the technology's emergence year. This approach yields a very similar allocation, as can be seen from comparing the figures in Panels A and B.

host a remarkable 40.2% of our pioneer locations.<sup>20</sup> Another cluster of major cities along the northeast corridor, New York Boston and Washington (DC), accounts for an additional 24.58%.

The geographic distribution of patenting related to our 29 disruptive technologies is even more skewed than that of general patenting, which, as discussed by Moretti (2019), is unevenly distributed geographically. Appendix Figure 5 depicts the population-normalized share for the top 20 CBSAs of patents linked to disruptive technologies, and the population-normalized share for all patents over the same period.

These differences can also be shown through summary statistics. The coefficient of variation of the geographic distribution of overall cite-weighted patenting is 1.21, while that of patents exposed to our 29 disruptive technologies is 1.42. Similarly, for overall patenting, it takes 12 CBSAs to account for 50% of all patents, while the top five urban regions produce 33.8% of all patents. By contrast, it takes only 7 CBSAs to account for 50% of all disruptive patents, and the top five urban regions alone represent 42.2%. When we look at the 189 technology-pioneer location pairs discussed above, the corresponding numbers are 5 and 54.5%.

Figure 5 Panels B through E continue to mark pioneer locations with hollow blue circles, but now also add the location of technology job postings in the start year of the technology (the average  $Normalized\ share_{i,\tau,0}$  across technologies at  $t = 0$ ), where darker dots correspond to a higher normalized share of jobs.<sup>21</sup> The figure shows a remarkable alignment between innovation and 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 developed.

The remaining panels (C-E) show the evolution of this relationship as the technology matures (in years 1-2, 3-4, and 5-6, respectively). Although pioneer locations retain a higher share of technology employment throughout this period, we see a gradual diffusion of technology job postings, away from the pioneer locations and spreading out across the country.

In Table 5, we explore this relationship more formally using the specification:

---

<sup>20</sup> See Figure 6 for the percentage of pioneer-technology pairs by location. Also note that, despite California's exceptional role, all of our main results are robust to removing California from the sample.

<sup>21</sup> To facilitate comparison between panels, we calculate this average of normalized shares only for the 13 technologies that emerge during our Burning Glass sample and for which we have at least six years of data, that is, those emerging between 2007 and 2014.

$$Normalized\ share_{i,\tau,t} = \beta_1 Pioneer_{i,\tau} + \beta_2 Pioneer_{i,\tau}(t - t_{0,\tau}) + \delta_i + \delta_\tau + \delta_t + \varepsilon_{i,\tau,t} \quad (4)$$

where  $i$  denotes a CBSA,  $\tau$  denotes one of our 29 technologies,  $t$  denotes year, and  $t_0$  denotes year of emergence for the technology.  $Pioneer_{i,\tau}$  is a dummy which denotes the pioneer status of a CBSA-technology pair. In all specifications in Table 5, we control for technology, CBSA, and year fixed effects.

In Column 1, we see that while there is diffusion over time, the initial CBSAs where the new technology was invented retain their privileged positions. More specifically, the  $Normalized\ share_{i,\tau,t}$  of a technology's job postings is about 92 percentage points higher in its pioneer locations on average throughout the lifecycle of the technology. Table 5, Column 2, however, shows that the initial advantage of pioneer locations in job postings (231 percentage points at the year of emergence) decreases significantly over time, at a rate of about 6% per year ( $\beta_2/\beta_1 = 0.063$ , s.e.=0.063). The initial advantage thus fully dissipates in about 15.8 years.

In columns 3 and 4, we find similar patterns for technology hiring near Pioneer locations. Column 3 shows that in CBSAs within 100 miles of a technology's pioneer location we find about 16 percentage points higher levels of job postings in that technology throughout the lifecycle of the technology – and this regional spillover decays in tandem with the pioneer advantage over time, at about 4% per year (column 4).

### *Skill Broadening*

We next turn to examining the skill component of technology job postings over time. Figure 7 plots a measure of skill requirements of these job postings (the red circles). We compute for each SOC code, as reported by Burning Glass, the corresponding skill level as reported in the U.S. Census Bureau's American Community Survey for 2015. When multiple SOC codes are associated with a given technology  $\tau$  in year  $t$ , we compute a weighted average of the skill measure as follows:

$$Skill_t^\tau = \frac{\sum_o N_{o;t}^\tau \chi_{o;2015}}{\sum_o N_{o;t}^\tau}$$

where  $o$  is a Census SOC code,  $N_{o;t}^\tau$  is the number of Burning Glass job postings exposed to technology  $\tau$  and SOC code  $o$  at time  $t$ , and  $\chi_{o;2015}$  is the average skill level for SOC  $o$ , as measured by the 2015 ACS sample. We consider four different measures of skill at the SOC level: the share of college educated persons (baseline), the share of persons with post-graduate

qualifications, the average wage of persons, and the average years of schooling for persons in the SOC.<sup>22</sup>

Figure 7 plots the percentage of college-educated persons associated with job postings against the year since emergence on a technology-by-technology basis. The figure suggests that for the vast majority of technologies, there is a sharp decline in the skill level required for the positions associated with new technologies over time. Even in cases where demand for positions is sharply accelerating (such as AI and virtual reality), the share of skilled positions subsides over time. These results are consistent with the view that new technologies typically start with high-skill occupations and then involve larger parts of the workforce over time. The figure also shows a few notable exceptions to this general pattern: positions exposed to the Online Streaming, Cloud Computing, Search Engine, and Software Defined Radio technologies show no evidence of a declining average skill level over time (in fact, the trend for Online Streaming appears significantly positive).

We summarize this information by presenting a binned scatterplot in Figure 8. This depiction shows the relationship across all 29 technologies between time elapsed after the emergence year and the mean share of the postings for college-educated persons. It shows, on average, a strong negative linear trend, implying a declining requirement for a college-trained workforce as technologies mature.

Table 6 looks at this relationship formally. The sample consists of annual observations of each technology between 2007 and 2019. Here, we use the alternative measures of the skills required in the job postings associated with a given technology: the dependent variables include the share of the weighted SOC classes that are college educated (as in the figures above), the share with graduate degrees, mean wages, and the mean years of schooling. Each regression uses as the key independent variable the years since the emergence date and controls for technology and calendar year fixed effects. The specification again follows Table 4 regarding the criteria for inclusion in the analysis and weighting.

Using each measure, there is a strong negative relationship between the maturity of the technology and the reliance on a highly educated workforce. For instance, Columns 1 and 3 show that each

---

<sup>22</sup> The BG data also includes an indicator for college requirement for a subset of observations. However, since this subset is quite limited we prefer using SOC codes to generate this variable.

additional year since the emergence of the technology is associated with a fall of about 0.96 percentage points in the share of job postings requiring a college education (an annual decline of -1.71%) and a decline of \$1,023 in annual wages (measured in 2015 constant dollars) for the job postings associated with the technology. Similarly, the share of job postings in occupations requiring a post-graduate degree declines by a rate of 1.80% per year on average.

This skill-broadening effect sheds an interesting light on how high-skilled labor is complementary with low-skilled work. While there is an important body of work highlighting the way in which technological change has favored high-skilled occupations and contributed to wage inequality (Acemoglu, 2002; Goldin and Katz, 2009; Acemoglu and Autor, 2011 are examples), the way in which the hiring associated with new technologies can transition over time highlights the dynamics in this relationship.<sup>23</sup>

#### *Differential Region-broadening by Skill Level*

We next explore the heterogeneity of our region-broadening and pioneer persistence facts across skill categories. We use the SOC codes to divide our sample of job postings into three categories using the share of college-educated persons in each SOC code. (Again, using information from the 2015 ACS.) We term these high (job postings for occupations with at least 60% college educated), medium (with 30% to 59% college educated), and low skilled (less than 30% college educated). For instance, almost all optometrists in the ACS are college educated: thus, all job postings for optometrists are allocated to the high-skill category. We then examine how the decline in the coefficient of variation described above changes after the emergence year, and how these shifts differ across skill levels.

Figure 9 takes a first look at these patterns. It again is a binned scatterplot of the coefficient of variation by year, but with the two extremes (low and high skill) of this three-fold division. It shows that the decline in the coefficient of variation across regions is substantially steeper for low-skilled jobs than that for high-skilled ones. While the low-skilled job postings rapidly disperse across the country, the higher-end ones remain more bunched together.

---

<sup>23</sup> In particular, our findings provide support for key assumptions in the literature on automation – that high-skill workers have a comparative advantage in new tasks, and that this advantage erodes as technologies mature (Acemoglu and Restrepo, 2020).



Table 7 studies these patterns in more detail. It emulates the structure of the specification in Table 4, but now breaks the observations of technologies into high and low-skill buckets (omitting the medium-skill bucket) and adding an interaction between the years since emergence variable and a dummy for low-skill occupations. All specifications show a significantly larger decline in concentration for lower-skill occupations. In terms of magnitudes, the annual decline in the coefficient of variation for low-skill job postings is more than three times larger than that for high-skill jobs, declining by 3.7% annually for low-skill jobs and only by 1.1% for high-skill jobs. Appendix Table 17 shows this specification separately for job postings in the three skill buckets. Again, high-skill professions show a less steep decline in geographic concentration, although the coefficient of variation declines significantly for all three groupings over time.

We obtain similar results for the persistence of pioneer advantage result in Table 8. This table repeats the analysis of Table 5, column 2 separately for each bucket (low, medium, and high skill) of job postings. Rather than looking at dispersion, however, it focuses on the related concept of persistence of the pioneer region. Consistent with the earlier results, we find the decline in initial pioneer advantage is greater in the case of low-skilled than in high-skilled positions. The degradation in geographic concentration is about 6.7% per year for low-skill job postings, which is about twice the magnitude for high-skill job postings (3.5%). That is, pioneer locations where disruptive technologies were developed retain a long-term advantage in attracting job postings in that technology, particularly in high-skill occupations. The estimates in column 3 suggest this high-skill advantage dissipates fully only after 28.6 years.

### *Properties of Pioneer Locations*

Before turning to the diffusion of disruptive technologies in other dimensions, we explore the characteristics of pioneer locations where disruptive technologies were developed (and also did the bulk of their hiring at the time of the emergence date). In particular, we highlight that there is a strong relationship between academic centers and the pioneer locations where nascent disruptive technologies originate.

To this end, we calculate for each CBSA-technology pair the number of patents exposed to that technology ten years prior to the technology's emergence year. (Recall our definition of pioneer locations is based on this variable: a dummy that is one for locations that account for 50% of a technology's patents in that year). We normalize this number by CBSA population in the

emergence year and then regress this ratio (patents in technology per 1000 inhabitants) on region characteristics in 2015 (using data from the ACS).

The key independent variables, which measure the presence of research universities and skilled persons in a CBSA, are the logarithm of the volume of university assets (standardized by population), the university enrollment (standardized by population), the share of the population in the CBSA that is college educated or has a post-graduate degree, and the log average wage in the CBSA.<sup>24</sup> All specifications control for technology-specific fixed effects.

Panel A of Table 9 shows a strong cross-sectional pattern. Regions with a greater academic or skill presence—whether manifested by greater research university presence or a more educated workforce—were more likely to be involved in the early development of disruptive technologies. These patterns are illustrated graphically in Appendix Figure 6.

Perhaps more importantly, and consistent with our results above, Panel B shows that these same variables also account for higher per capita technology job postings in the emergence year. That is, the same variables that account for the location of innovative activity also account for early employment in that technology.

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

In this section, we characterize the spread of disruptive technologies across industries, occupations, and firms. First, we compare the region-broadening result against broadening across industries, occupations, and firms; second, similar to Table 5, we also study initial advantage of pioneers, separately defined across the four segments, and the degradation in this advantage over time.

To that end, we extend the definition of  $Normalized\ share_{i,\tau,t}$  in Section 3 to NAICS four-digit industries, SOC six-digit occupations, and firms for each technology ( $\tau$ ) and time ( $t$ ). While the former two variables are included in the BG data (in each case, we use the finest level of

---

<sup>24</sup> We obtain university data for 642 research universities from the U.S. National Science Foundation’s Higher Education Research and Development Survey (HERD) and from the Integrated Postsecondary Education Data System (IPEDS) surveys provided by the U.S. Department of Education’s National Center for Education Statistics (NCES), and map these universities to CBSAs. Research universities are defined as “public and private nonprofit postsecondary institutions in the United States, Guam, Puerto Rico, and the U.S. Virgin Islands that granted a bachelor’s degree or higher in any field; expended at least \$150,000 in separately budgeted R&D in FY 2015; and were geographically separate campuses headed by a president, chancellor, or equivalent.” We normalize university assets and the university enrollment by CBSA population from the ACS at the time of the year of emergence. We obtain skill level variables for a particular CBSA from the ACS, by normalizing the share of graduate and post graduate persons in a CBSA by the total number of persons in the CBSA. For further details, refer to Appendix 3.2.

disaggregation available from BG), the latter relies on our own matching algorithm described in Section 2.

We then measure the coefficient of variation of *Normalized share* <sub>$i,\tau,t$</sub>  across the segments. Because the number of firms posting job advertisements online expands over time (with more and more small firms appearing in the BG data over time), we stratify our firm-technology-year sample by including only firms that post at least one job in each of our sample-years, before calculating the coefficient of variation.<sup>25</sup> This step focuses attention on 10,231 larger firms which on average post 1,628 jobs per year, effectively excluding variation coming from small and medium-sized businesses.

Table 10, Panel A shows the results of a regression of the coefficient of variation calculated for each technology ( $\tau$ ) and time ( $t$ ) on year since emergence. Column 1 shows our already established results for locations for comparison.<sup>26</sup> We find that while there is a decline in concentration as measured by coefficient of variation for all four segments, there appears to be a larger decline across locations and firms (Columns 1 and 4) than across industries and occupations (Columns 2 and 3). While the coefficient of variation declines on average by 2.48% and 2.32% for CBSAs and firms, respectively, the corresponding declines are 1.06% and 0.81% for (four-digit NAICS) industries and (six-digit SOC) occupations, respectively. Figure 10 illustrates these patterns graphically, and Appendix Figures 7 through 9 illustrate them technology by technology.

While it is perhaps natural to expect disruptive technologies to spread faster across firms and space than they do across industries and occupations, any quantitative comparison of course depends on the classifications of industries and occupations used. Appendix Figures 7 through 9 show some differences across technologies in their diffusion across industries and occupations. For example, the 3D Printing, Computer Vision, and Wi-Fi technologies show a clear decrease in concentration across industries over time.

In Table 10, Panel B, we estimate specification (1) for all four dimensions to examine the initial hiring advantage of pioneer cells in the four segments. The pioneer cells, as defined before, are

---

<sup>25</sup> Hershbein and Kahn (2018) discuss this fact in some detail. The general increase in coverage of the BG data over time should not affect any of our main results. We discuss robustness to various weighting schemes in detail in Section 5.

<sup>26</sup> In order to avoid calculating coefficients of variation for unreasonably sparse data, we only keep technology x year observations with at least 100 postings with industry coverage. This issue arises because BG provides NAICS codes for only 41% of all postings, as noted above.

ones that excel in initial technology-related inventions. More specifically, we define pioneer cells (occupations, industries, and firms) as those which collectively accounted for 50% of the patent grants associated with a given technology in the ten years before its emergence.

To determine the pioneer cells, we merge various public-use datasets to assign patents to our three additional segments of industries, occupations and firms: for industries, we allocate patents to individual NAICS four-digit industries by mapping patents to Compustat firms (since patents themselves do not contain industry codes), and then from firms to industries. To obtain the patent-to-Compustat match, we use the crosswalk provided by Autor et al. (2020), who use the Bing search engine to match assignee names from patents to Compustat firms. A total of 44% of all patents exposed to any one of our 29 technologies match to a Compustat firm, so that this procedure implicitly assumes that the distribution of patents across industries is similar for Compustat firms as for all firms.<sup>27</sup> Once patents are matched to firms, we then link to industries using the Compustat Segments dataset, which gives firms' breakdown of sales across NAICS four-digit industries. So, for example, if a patent is owned by "Apple North America," it is matched by Bing to "Apple Inc.," and then allocated proportionally to Apple's NAICS four-digit industries by its sales breakdowns (83% to "Computer and peripheral equipment manufacturing" and 17% to "Electronics and Appliance Stores").

For occupations, we further construct an industry-to-occupation crosswalk from employment data within an occupation-industry cell from the Occupational Employment Statistics. We assume that the share of patenting in an industry allocated to an occupation is the same as the share of employment allocated to an occupation. We thus can, calculate the share of patents for a particular technology allocated to an occupation<sup>28</sup>.

Finally, for firms, we string match patent assignees from USPTO to firm names in job postings. (See Appendix 1.4 for details.) Using this procedure, we are able to match 36% of all patents assigned to U.S. inventors between 1976 and 2016 to 30,123 unique firms in our sample.

---

<sup>27</sup> In order to compare patents by Compustat and non-Compustat firms, we analyze the share of patents by Compustat-firms across technology classes. We find that for the median technology class, about 50% of patents are produced by Compustat firms, and that the distribution is quite homogenous: the 25<sup>th</sup> percentile is 39.0% and 75<sup>th</sup> percentile is 58.8%.

<sup>28</sup> We reweigh technology jobs in an occupation to match hiring in the U.S. economy for each two-digit occupation. Hiring in a two-digit occupation in the US economy is calculated using hiring in an industry in the Longitudinal Employer-Household Dynamics Census survey and then constructing a crosswalk between industry employment and occupation employment. For more details refer to Appendix 3.3.

Following our procedure for pioneer locations, we define pioneer industries, occupations, and firms for each technology as those with the most assigned patents in the ten years prior to the technology’s emergence year that collectively account for 50% of the matched patents in a given disruptive technology. Appendix Table 18 shows the top pioneer industry and occupation for each technology. For example, the top pioneer industry for 3D Printing is “Computer and Peripheral Equipment Manufacturing” (accounting for 41.9% of early patents) and its top occupation is “Mechanical Engineers,” while that of Fracking is “Oil and Gas Extraction” (accounting for 88.1% of early patents) and “Geoscientists,” respectively.

The analysis in Table 10, Panel B shows that pioneering cells have a strong initial advantage in job postings for all four segments. Over time, these again degrade significantly. The estimated degradation in the advantage ( $\beta_2/\beta_1$  from equation (4)) for locations is 6.2%, compared to 4.4% for firms, 4.0% for industries, and 3.4% for occupations.

Taken together, this evidence suggests disruptive technologies initially generate hiring that is highly localized by location, firm, industry and occupation. Over time, this hiring disperses, particularly across locations and across firms.

#### *Firm Rehoming towards Pioneer Locations*

As a final analysis, we explore one of the mechanisms behind the region-broadening results: the rehoming of firms towards pioneer locations using a case study.

More specifically, we consider the geographic footprint of Ford Motor Company and General Motor Corporation before and after the emergence year of the autonomous cars technology (2014). In Figure 11, we plot these firms’ job postings in three groups of places: (a) the three autonomous car pioneer locations, San Jose (CA), San Francisco (CA), and Boston (MA) (but excluding Detroit (MI)); (b) their headquarters, Detroit (MI), and (c) all other locations. Postings in red are before the emergence year of autonomous car technology, and postings in blue are post-emergence year. Black crosses in the picture denote the share of job postings exposed to autonomous cars post emergence year.

The figure shows that both firms, traditionally concentrated in Detroit, shifted their geographic footprint towards the autonomous cars pioneer locations, particularly in Silicon Valley (San Jose and San Francisco). A large share of new job postings in the pioneer locations involved

autonomous car technologies, accounting for 22% and 65% of Ford and GM postings respectively (compared to less than 5% in all other locations). The data thus suggest that the purpose of both firms' expanding presence in autonomous cars' pioneer locations related to this new technology.

## 5. Additional Robustness Checks

Before concluding, we perform a number of additional robustness checks for our primary results: “region broadening,” “pioneer-location persistence,” “skill broadening,” and “differential region-broadening by skill level.”

First, we replicate our results using our “unsupervised” approach to defining technologies. That is, we treat each of the original 305 technical bigrams we obtained from our algorithm intersecting the texts of patents and earnings calls as a separate technology, without attempting to group or otherwise audit these bigrams. The goal of this exercise is to replicate our main findings in a dataset created without any human intervention.

In Table 11, columns 1 through 4 of Panel A replicate the main specifications of Tables 4 through 8, respectively. We find that all the coefficients of interest are qualitatively and quantitatively similar to our main specification. Column 1 shows our region-broadening result, regressing each technology-year's coefficient of variation across locations on the number of years since the emergence of the technology. The estimated coefficient ( $-0.140$ ,  $s.e.=0.017$ ) implies a 2.54% reduction in concentration in technology job postings per year, compared to 2.21% in our baseline specification (Table 4, column 1). Similarly, the estimates in column 2 imply a large advantage of pioneer locations in job postings that decreases at a rate of 6.0% per year, compared to 6.6% in Table 5, column 3. Column 3 also shows significant skill broadening over time, with a decreasing share of job postings that require a college education over the life-cycle of the technology. However, the estimate here ( $-0.325$ ,  $s.e.=0.099$ ) is only one third the size of that in Table 6, column 1. Finally, column 4 shows that the geographic concentration of low-skill jobs exposed to disruptive technologies decays significantly faster than that of high-skill jobs, though the coefficient of interest is again somewhat smaller ( $-0.108$ ,  $s.e.=0.028$  vs.  $-0.167$ ,  $s.e.=0.048$  in our baseline specification).

Panel B of Table 11 replicates the results of Table 10, estimating the spread of disruptive technologies across industries, occupations, and firms. The results are again similar, although this

unsupervised approach yields somewhat faster spread across occupations than in our baseline specification.

We conclude that the human judgement that we exerted to enable us to measure the spread of specific technologies has no bearing on the validity of our main stylized facts about disruptive technologies as a whole.

Second, in Appendix Table 19, we check for robustness with respect to our methodology for calculating the year of emergence for our technologies, with respect to the missing years (2008 and 2009) in the BG sample, and with respect to standard errors:

In Panel A, we find that our results are robust to calculating years of emergence exclusively from patents instead of earnings calls. To calculate this alternative measure, we use our patent data, which extends back to 1975. The year of emergence for each technology is here calculated as the year in which the share of U.S. patents exposed to the technology reached 50% of their maximum value between 1976 and 2015.

In Panel B, we find that our results are robust to excluding 2007, the first year of availability of Burning Glass job postings and immediately before the missing BG job postings in 2008 and 2009.

In Panel C, we check for robustness of standard errors and find that if anything the statistical significance is stronger with robust standard errors (vs. clustered standard errors in the baseline specification).

Third, we deal with the potential concern that some of our analyses may reflect changes in the composition of the job announcements in Burning Glass, not hiring overall. Appendix Figure 10a shows that the number of job postings in the BG dataset began increasing sharply in the mid-2010s (the blue line), which could reflect an increase in the share of jobs posted online. We note, however, that this trend also parallels the increase in overall U.S. job openings after the 2008-09 recession, as reported by the Job Openings and Labor Turnover Survey (the red line).

A more substantive compositional concern is raised by Appendix Figure 10b. The figure shows that much of the growth in Burning Glass online job postings was driven by job postings in low-skill occupations. It is natural to speculate that many of these jobs may have previously not been posted online. Thus, the increase in BG postings shown in Appendix Figure 10a may reflect both increasing overall U.S. hiring and a growing tendency for lower-skill job announcements to be

posted online. It is thus natural to wonder whether the changing composition of BG job announcements may have impacted the results above.

After three additional robustness checks, we do not believe these changes affect the results in our analyses. First, it is important to note, as demonstrated in Appendix Figure 10c, that the compositional patterns documented in Appendix Figure 10b are much less pronounced among job announcements associated with our 29 technologies. Second, our entire analysis uses the normalized share of job postings (except skill broadening), and controls throughout for year fixed effects. The normalization and year controls should address many of these compositional concerns. As a final check for our skill broadening result, we reweight the occupations in our sample to match hiring in that occupation in the U.S. economy. We compute hiring in each occupation by using hiring in each industry in the Longitudinal Employer-Household Dynamics Census survey and then constructing a crosswalk between industry employment and occupation employment using the Occupational Employment Statistics Census survey. In Appendix Table 20, we find that our skill-broadening results are robust to this reweighting exercise.

We also check our primary results for sensitivity with respect to our “rising” cut off in earnings calls. To get to our list of 305 unsupervised bigrams, we keep bigrams which appear at least ten times as frequently in their peak year as in the first year of the earnings call data in 2002. In Appendix Table 21, we vary this to keep bigrams which appear at least 100 times, 20 times, 10 times, 6 times and 5 times as frequently in their peak year as in the first year of the earnings call data in 2002. In Appendix Table 22, we find that our primary results are fully robust to changes in varying these cut-offs.

As a final check of our broadening results, we check their sensitivity to technology selection: in other words, could the results be driven by a handful of industries out of our 29? To do this, we exclude three technologies at a time and recalculate the degradation in coefficient of variation, this provides us with 7,308 permutation estimates. In Appendix Figure 11, we plot the 5<sup>th</sup> and 95<sup>th</sup> percentile of these jackknife estimates, and show that the results are robust to randomly removing a subset of technologies.

## **6. 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 millions of patents, job postings, and earnings conference calls over the past two decades, we introduce in this paper an approach to understand which new technologies affect businesses and to trace their diffusion across regions, industries, occupations, and firms. We can then map the spread of disruptive technologies in these dimensions, focusing on the hiring associated with each important innovation.

We highlight five main conclusions. First, the locations where disruptive technologies are developed are geographically highly concentrated, with a handful of urban areas contributing the bulk of the early patenting and early employment within each technology. Second, despite this initial concentration, jobs relating to use or production of the new technologies gradually spread out geographically. Third, while initial jobs associated with a given technology are typically high-skilled, over time the mean required skill levels of the new jobs declines. Fourth, these trends towards region and skill broadening are related: low-skill jobs associated with a given technology spread out geographically significantly faster than high-skill ones. Finally, because of the slower spread of high-skill jobs, disruptive technologies continue to offer long-lasting benefits for their pioneer locations, which retain a long-term advantage in these high-quality jobs for multiple decades.

Beyond these core results of our analysis, the development and spread of disruptive technologies are key objects of interest in multiple fields of economics. We therefore hope that the data we provide as part of this paper may prove useful to address a range of additional research questions in the study of economic growth, inequality, entrepreneurship, and firm dynamics.

## References

- Acemoglu, Daron. "Directed technical change." *Review of Economic Studies* 69 (2002): 781-809.
- 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, pages 1043-1171 (2011).
- Acemoglu, Daron, and Joshua Linn. "Market size in innovation: Theory and evidence from the pharmaceutical industry." *Quarterly Journal of Economics* 119 (2004): 1049-1090.
- Acemoglu, Daron, and Pascual Restrepo. "Robots and jobs: Evidence from US labor markets." *Journal of Political Economy* 128 (2020) 2188-2244.
- Aghion, Philippe, Ufuk Akcigit, Antonin Bergeaud, Richard Blundell, and David Hemous. "Innovation and Top Income Inequality." *Review of Economic Studies* 86 (2019): 1–45.
- Agrawal, Ajay. Joshua Gans, and Avi Goldfarb, editors. *The Economics of Artificial Intelligence: An Agenda*. Chicago: University of Chicago Press (2019).
- Akerman, Anders, Ingvil Gaarder, and Magne Mogstad. "The skill complementarity of broadband internet." *Quarterly Journal of Economics* 130 (2015): 1781-1824.
- Arora, Ashish, Sharon Belenzon and Lia Sheer. "Knowledge spillovers and corporate investment in scientific research." *American Economic Review* 111 (2021) 871-898.
- Atalay, Enghin, Phai Phongthientham, Sebastian Sotelo, and Daniel Tannenbaum. "The evolution of work in the United States." *American Economic Journal: Applied Economics* 12 (2020): 1-34.
- Autor, David H., David Dorn, Gordon H. Hanson, Gary Pisano, and Pian Shu. "Foreign competition and domestic innovation: Evidence from US patents." *American Economic Review: Insights* 2 (2020): 357-74.
- 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." *The quarterly journal of economics* 131.4 (2016): 1593-1636.
- Bekkerman, Ron, and James Allan. Using bigrams in text categorization. Technical Report IR-408, Center of Intelligent Information Retrieval, UMass 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.
- Bessen, James, Iain Cockburn and Jennifer Hunt, "Is Distance from Innovation a Barrier to Adoption of Artificial Intelligence", BU mimeo, 2021.
- 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, L., Kelly, B. T., Manela, A., & Xiu, D. (2020). The structure of economic news (No. w26648). National Bureau of Economic Research.
- 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.
- 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.
- Glaeser, Edward L., Sari Pekkala Kerr, and William R. Kerr. "Entrepreneurship and urban growth: An empirical assessment with historical mines." *Review of Economics and Statistics* 97 (2015): 498-520.
- Goldin, Claudia D., and Lawrence F. Katz. "The origins of technology-skill complementarity," *Quarterly Journal of Economics* 113 (1998): 683-732
- Goldin, Claudia D., and Lawrence F. Katz. *The Race Between Education and Technology*. Cambridge, Harvard University Press (2009).
- Gompers, Paul, Josh Lerner, and David Scharfstein. "Entrepreneurial spawning: Public corporations and the genesis of new ventures, 1986 to 1999." *Journal of Finance* 60 (2005): 577-614.
- Gordon, Robert J. *The Rise and Fall of American Growth: The U.S. Standard of Living since the Civil War*. Princeton, Princeton University Press (2016).
- Greenstone, Michael, Richard Hornbeck and Enrico Moretti, "Identifying agglomeration spillovers: Evidence from winners and losers of large plant openings." *Journal of Political Economy*, 118 (2010): 536-598.
- Griliches, Zvi. "Hybrid corn: An exploration in the economics of technological change." *Econometrica* 25 (1957): 501-522.
- Hall, Bronwyn H. "Innovation and diffusion." Working paper no. 10212, National Bureau of Economic Research (2004).
- Handley, K., & Li, J. F. (2020). Measuring the effects of firm uncertainty on economic activity: New evidence from one million documents (No. w27896). National Bureau of Economic Research.
- Hansen, S., Ramdas, T., Sadun, R., & Fuller, J. (2021). The Demand for Executive Skills (No. w28959). National Bureau of Economic Research.

- 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 Alexander, 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.
- Jaffe, Adam B., Manuel Trajtenberg, and Rebecca Henderson. "Geographic localization of knowledge spillovers as evidenced by patent citations." *Quarterly Journal of Economics* 108 (1993): 577-98
- 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* (forthcoming).
- Krueger, Alan B. "How computers have changed the wage structure: Evidence from microdata, 1984-1989." *Quarterly Journal of Economics* 108 (1993): 33-60.
- Lanjouw, Jean O., Ariel Pakes, 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 corporate finance and beyond." *Review of Financial Studies* (forthcoming).
- Machin, Stephen, and John van Reenen. "Technology and changes in skill structure: Evidence from seven OECD countries." *Quarterly Journal of Economics*, 113 (1998) 1215-44.
- 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.

- Michaels, Guy, Natraj, Ashwini and Van Reenen, John. "Has ICT polarized skill demand? Evidence from eleven countries over 25 years." *Review of Economics and Statistics* 96 (2014): 60-77
- Mikolov, Tomas, Kai Chen, Greg Corrado, and Jeffrey Dean. "Efficient estimation of word representations in vector space." *arXiv preprint arXiv:1301.3781* (2013).
- Moretti, Enrico. "The effect of high-tech clusters on the productivity of top inventors." Working paper no. 26270, National Bureau of Economic Research (2019).
- Moscona, Jacob. "Environmental catastrophe and the direction of invention: Evidence from the American Dust Bowl." Unpublished working paper, Massachusetts Institute of Technology (2020).
- Mokyr, Joel. *The Lever of Riches: Technological Creativity and Economic Progress*. New York: Oxford University Press (1992).
- Moser, Petra, Alessandra Voena, and Fabian Waldinger. "German Jewish Émigrés and US Invention." *American Economic Review* 104 (2014): 3222–55.
- 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.
- Sautner, Z., van Lent, L., Vilkov, G., & Zhang, R. (2020). Firm-level climate change exposure.
- 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): 1825-1883.
- 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.
- United States Patent and Trademark Office. "Performance and Accountability Report." (2020).

Abis, Simona, and Laura Veldkamp. "The changing economics of knowledge production."  
Available at SSRN 3570130 (2020).

Webb, Michael. "The Impact of Artificial Intelligence on the Labor Market." Unpublished  
working paper, Stanford University (2020).

## Tables and Figures

Table 1 – Most frequent technical bigrams in earnings calls

Bigram	# ECs	Technology group
mobile devices	6597	Smart devices
machine learning	2860	Machine learning/AI
cloud computing	2781	Cloud computing
cloud services	2450	Cloud computing
quality metrics	2029	NA
flow profile	1966	NA
smart phones	1957	Smart devices
mobile platform	1605	Smart devices
public cloud	1569	Cloud computing
social networking	1548	Social networks
smart grid	1441	NA
cloud service	1393	Cloud computing
connected devices	1304	Smart devices
cloud infrastructure	1136	Cloud computing
carbon footprint	1071	NA
nand flash	1002	NA
virtual reality	903	Virtual reality
digital channel	896	NA
delivery network	887	NA
social networks	883	Social networks
autonomous driving	839	Autonomous cars
smart devices	765	Smart devices
active user	735	Social networks
augmented reality	730	Virtual reality
mobile payment	717	Mobile payment
cloud environment	668	Cloud computing
production site	664	NA
ethanol production	662	NA
power outage	643	NA
multiple segments	595	NA

Notes: This table lists the top 30 most frequent technical bigrams in Earnings Calls (ECs) that are associated with disruptive technologies. Column 2 gives the number of ECs each bigram is mentioned in. Column 3 reports each bigrams's technology grouping under the supervised approach described in section 2.3. Those labeled with 'NA' are excluded from the supervised but not the unsupervised approach.



Table 2 – Technologies by number of job postings

Technology	Postings
Cloud computing	3,684,901
Social networking	3,457,390
Smart devices	2,376,510
Machine learning/AI	679,776
Search engine	535,784
Online streaming	487,731
Wi-Fi	388,844
Electronic gaming	247,201
Solar power	201,296
Injection molding	190,538
Hybrid vehicle/Electric car	118,550
Touch screen	109,538
RFID	80,894
Computer vision	76,350
GPS	65,922
Mobile payment	65,482
Virtual reality	61,102
3D printing	57,904
Autonomous cars	52,974
Lane departure warning	32,107
Lithium battery	16,926
Software defined radio	14,187
Drug conjugates	10,603
Fracking	8,966
Millimeter wave	6,161
OLED display	5,528
Bispecific monoclonal antibody	2,702
Inkjet printing	2,583
Wireless charging	1,649
Stent graft	1,270
Fingerprint sensor	711

Notes: This table lists our 29 technologies (in Column 1) and the number of job postings from Burning Glass 2007-2019 that mention the technology (in Column 2).

Table 3 – Technology keywords

Technology	Keywords
3D printing	3d printer; 3d printing; additive manufacturing; 3d printed
Autonomous cars	Self-driving car; robot car; autonomous vehicles; autonomous car; autonomous cars; automated driving; driverless car; autonomous driving; autonomous vehicle; driverless truck
Bispecific monoclonal antibody	bispecific monoclonal; the bispecific; bispecific antibody
Cloud computing	paas; cloud infrastructure; distributed cloud; cloud provider; cloud offerings; cloud service; cloud applications; community cloud; private cloud; public cloud; cloud deployments; cloud environments; cloud management; cloud services; cloud security; enterprise class; iaas; hybrid cloud; cloud platform; cloud providers; cloud hosting; personal cloud; enterprise network; cloud computing; cloud based; saas; cloud storage; enterprise applications; cloud solution; enterprise cloud; cloud solutions; cloud deployment
Computer vision	pose estimation; motion estimation; visual servoing; facial recognition; gesture recognition; computer vision; image recognition; sensor fusion; object recognition
Drug conjugates	kinase inhibitor; drug conjugate; antibody drug; drug conjugates
Electronic gaming	social game; video games; social games; video game; game content; electronic gaming; gaming products
Millimeter wave	millimeter wave
Fingerprint sensor	fingerprint sensor; fingerprint scanner
Fracking	fracking; fraccing; hydrofracking; hydrofracturing; hydraulic fracturing;

Notes: The table lists 10 of our technologies in alphabetical order (in Column 1) and the bigrams used to identify them in text of earnings calls, patents, and job postings (in Column 2). Any single words (unigrams) given denote all possible bigrams containing that unigram.

Table 4 - Region broadening

	(1)	(2)	(3)
	Coefficient of Variation	$\frac{Mean(NS)_{Top\ 5}}{Mean(NS)_{All}}$	Share CBSAs with ( $NS < 1\%$ )
<i>Years since emergence</i> <sub><math>\tau, t</math></sub>	--0.105*** (0.027)	--1.078*** (0.338)	--0.028*** (0.006)
R2	0.861	0.776	0.927
N	287	287	287
Tech FE	YES	YES	YES
Year FE	YES	YES	YES
Mean	4.74	53.33	0.67
% Mean per year	2.21%	2.02%	4.18%

Notes: This table reports results from a regression of three separate measures of geographic concentration, calculated over  $Normalized\ share_{i,\tau,t}(NS) = \frac{share\ jobs\ exposed_{i,\tau,t}}{share\ jobs\ exposed_{\tau,t}}$ , where  $i$  denotes a location (CBSA),  $\tau$  technology, and  $t$  time. The three measures are: coefficient of variation, the mean normalized share of technology postings in the top 5 CBSAs relative to the mean normalized share across all CBSAs, and the share of CBSAs with a (negligible) normalized share of technology job postings of less than 1%. Observations are weighted by the square root of total technology job postings in a given year. The normalized share is capped at 99<sup>th</sup> percentile of non-zero observations. Standard errors are clustered by technology. The last row specifies the magnitude of the coefficient of *Years since emergence* <sub>$\tau, t$</sub>  as a percentage of the sample mean per year.

Table 5 – Persistent advantage of Pioneer locations

	(1)	(2)	(3)	(4)
	Normalized Share			
<i>Pioneer</i> <sub><math>i,\tau</math></sub>	0.918*** (0.285)	2.313*** (0.580)	0.938*** (0.285)	2.340*** (0.581)
<i>Pioneer</i> <sub><math>i,\tau</math></sub> * <i>Years since emergence</i> <sub><math>\tau, t</math></sub>		--0.146*** (0.042)		--0.147*** (0.042)
( $\leq 100$ miles to <i>Pioneer</i> <sub><math>i,\tau</math></sub> )			0.155*** (0.038)	0.260*** (0.056)
( $\leq 100$ miles to <i>Pioneer</i> <sub><math>i,\tau</math></sub> ) * <i>Years since emergence</i> <sub><math>\tau, t</math></sub>				--0.011** (0.005)
R2	0.074	0.075	0.075	0.076
N	266,467	266,467	266,467	266,467
$\beta(Pioneer_{i,\tau} * (Years\ since\ emergence_{\tau,t})) / \beta(Pioneer_{i,\tau})$		-0.063*** (0.007)		-0.063*** (0.008)
Tech FE	YES	YES	YES	YES
Year FE	YES	YES	YES	YES
CBSA FE	YES	YES	YES	YES

Notes: This table reports results from a regression of the *Normalized share* <sub>$i,\tau,t$</sub>  (for each CBSA, technology, and year) on Pioneer status of the CBSA and its interaction with the year since technology's emergence. Normalized share is capped at 99<sup>th</sup> percentile of non-zero observations. ( $\leq 100$  miles to *Pioneer* <sub>$i,\tau$</sub> ) is a dummy which takes value 1 for non-Pioneer CBSAs which are within 100 miles of any Pioneer location. Standard errors are clustered by CBSA. The expression  $\beta(Pioneer_{i,\tau} * (Years\ since\ emergence_{\tau,t})) / \beta(Pioneer_{i,\tau})$  is calculated by dividing the coefficient on *Pioneer* <sub>$i,\tau$</sub>  \* *Years since emergence* <sub>$\tau, t$</sub>  by the coefficient on *Pioneer* <sub>$i,\tau$</sub>  -- standard errors for this expression are calculated using the delta method.

Table 6 – Skill broadening

	(1)	(2)	(3)	(4)
	Share of college educated * 100	Share of post graduate * 100	Average wage	Average schooling
<i>Years since emergence</i> $_{\tau,t}$	-0.954*** (0.260)	-0.361*** (0.121)	-1,022.929*** (241.521)	-0.050*** (0.014)
R2	0.847	0.878	0.845	0.859
N	287	287	287	287
Tech FE	YES	YES	YES	YES
Year FE	YES	YES	YES	YES
Mean	55.90	19.95	64,463	15.07
%Mean/year	-1.71%	-1.80%	-1.59%	-.33%

Notes: This table reports the results from a panel regression of the skill composition of technology job postings in a given technology in a given year ( $Skill_t^T$ , see section 3), on the years since the emergence of the technology. The specification excludes all technology-year observations prior to the emergence year of the technology. Observations are weighted by the square root of the total number of technology job postings in a given year. All specifications control for technology and year fixed effects. Standard errors are clustered by technology.

Table 7 – Region Broadening – High vs. low skill

	(1)	(2)	(3)
	Coefficient of Variation	$\frac{Mean(NS)_{Top\ 5}}{Mean(NS)_{All}}$	Share CBSAs with ( $NS < 1\%$ )
$Years\ since\ emergence_{skill,\tau,t} * 1\{skill = low\}$	-0.167*** (0.048)	-2.218*** (0.621)	-0.022*** (0.006)
$Years\ since\ emergence_{skill,\tau,t}$	-0.074** (0.036)	-0.657 (0.464)	-0.017*** (0.005)
R2	0.773	0.653	0.827
N	567	567	567
Skill FE	YES	YES	YES
Tech FE	YES	YES	YES
Year FE	YES	YES	YES
Mean	6.53	16.85	0.78
% Mean/year	2.56%	13.16%	2.82%

Notes: This table reports results from regressions analogous to Table 4, but now breaks the observations of technology-years into high and low-skill buckets (omitting the medium-skill bucket) and adding an interaction between the years since emergence variable and a dummy for low-skill occupations. An observations is thus a skill-level-technology-year. The interaction with  $1\{skill = low\}$ , tests for differential concentration trends between low- and high-skill job postings. To calculate geographic concentration by skill x technology x year, we aggregate the job postings data over occupation, CBSA, and year, separately for high-skill occupations (with the share of college educated people  $> 60\%$ ) and low-skill occupations (with the share of college educated people  $< 30\%$ ). Finally, measures of concentration (as in table 4) are calculated over  $Normalized\ share_{i,\tau,t,skill} = \frac{share\ jobs\ exposed_{i,\tau,t,skill}}{share\ jobs\ exposed_{\tau,t,skill}}$  across CBSAs by skill group, technology, and time. All specifications exclude observations before the year of emergence of a technology. Observations are weighted by the square root of the total number of technology job postings in the skill-technology-year triple. All specifications control for skill-level, technology, and year fixed effects. Standard errors are clustered by technology.

Table 8 – Advantage of Pioneer locations by skill

	<i>Normalized share<sub>i,τ,t</sub></i>		
	(1) Low Skill	(2) Medium Skill	(3) High Skill
<i>Pioneer<sub>i,τ</sub></i>	1.607*** (0.403)	1.193*** (0.453)	1.108** (0.484)
<i>Pioneer<sub>i,τ</sub> * years since emergence<sub>τ,t</sub></i>	-0.108*** (0.030)	-0.057** (0.025)	-0.039* (0.020)
R2	0.053	0.044	0.049
N	181,598	181,598	181,598
$\beta(Pioneer_{i,\tau} * years\ since\ emergence_{\tau,t}) / \beta(Pioneer_{i,\tau})$	-0.067*** (0.007)	-0.048** (0.016)	-0.035** (0.014)

Notes: This table reports the results from regressions of *Normalized share<sub>i,τ,t,skill</sub>* on the Pioneer status dummy (and its interaction with time) for CBSA *i* and technology *τ*, separately for job postings in low skill (Column 1), medium skill (Column 2) and high skill (Column 3) occupations. To construct the sample at the skill-level x CBSA x year level, we aggregate the job postings data over occupation, CBSA, and year, separately for high-skill, medium-skill and low-skill occupations (see section 3 and the caption of Table 7 for details). All specifications exclude observations before the year of emergence of a technology and control for CBSA, technology and skill-level fixed effects. Standard errors are clustered by CBSA. The expression  $\beta(Pioneer_{i,\tau} * (Years\ since\ emergence_{\tau,t})) / \beta(Pioneer_{i,\tau})$  is calculated by dividing the coefficient on *Pioneer<sub>i,τ</sub> \* Years since emergence<sub>τ,t</sub>* by the coefficient on *Pioneer<sub>i,τ</sub>*. Standard errors for this expression are calculated using the delta method.

Table 9 – Initial patenting and job postings versus skill composition

Panel A:	(1)	(2)	(3)	(4)	(5)
	<i>Early patents per 1000 people<sub>i,t,0</sub></i>				
$\log(1 + \text{university assets (in \$1,000 per capita)})_i$	0.129*** (0.022)				
University enrollment per capita <sub>i</sub>		0.346*** (0.085)			
Share College Educated (in pct.) <sub>i</sub>			0.0178*** (0.0017)		
Share post graduate (in pct.) <sub>i</sub>				0.0421*** (0.0041)	
$\log(\text{wage}_i)$					1.004*** (0.117)
Observations	24,731	24,731	24,731	24,731	24,731
R-squared	0.107	0.093	0.158	0.162	0.133
Tech FE	YES	YES	YES	YES	YES
Panel B:	(1)	(2)	(3)	(4)	(5)
	<i>Early job postings per 1000 people<sub>i,t,0</sub></i>				
$\log(1 + \text{university assets (in \$1,000 per capita)})_i$	0.0595*** (0.0075)				
University enrollment per capita <sub>i</sub>		0.217*** (0.0313)			
Share College Educated (in pct.) <sub>i</sub>			0.0066*** (0.0006)		
Share post graduate (in pct.) <sub>i</sub>				0.0149*** (0.0015)	
$\log(\text{wage}_i)$					0.426*** (0.045)
Observations	24,759	24,759	24,759	24,759	24,759
R-squared	0.179	0.172	0.197	0.196	0.192
Tech FE	YES	YES	YES	YES	YES

Notes: The table presents results from a regression of initial patents per capita (in Panel A, *Early patents per 1000 people<sub>i,t,0</sub>*) and initial job postings per capita (in panel B, *Early job postings per 1000 people<sub>i,t,0</sub>*) in a CBSA (*i*) associated with a technology on values of various measures of skill and income for the CBSA. *Early patents per 1000 people<sub>i,t,0</sub>* are calculated by dividing patents applied for by inventors from a CBSA during the 10 years before the emergence of a technology by the population of the CBSA in the 2015 American Communities Survey. *Early job postings per 1000 people<sub>i,t,0</sub>* are calculated by dividing jobs posted in a CBSA during year of emergence of a technology by the population of the CBSA (again in 2015). University-related measures in Row 1 and Row 2 are calculated by aggregating university assets and enrollment over all universities in a CBSA, and share of college educated/post graduate in Row 3 and Row 4 are calculated as the share of people holding a college/postgraduate degree in a CBSA. Row 5 uses the log of average wages, where average wages for a CBSA are calculated as the average yearly income of a working person. The university data is from the U.S. National Science Foundation's Higher Education Research and Development Survey (HERD) and from the Integrated Postsecondary Education Data System (IPEDS) surveys provided by the U.S. Department of Education's National Center for Education Statistics (NCES). Income data is from American Communities Survey 2015. All specifications control for technology fixed effects. Standard errors are clustered by CBSA.

Table 10 – Broadening and Pioneer persistence: Comparison across different dimensions

Panel A:	<i>Coefficient of Variation<sub>τ,t</sub></i>			
	(1) Locations	(2) Industries	(3) Occupations	(4) Larger Firms
<i>Years since emergence<sub>τ,t</sub></i>	-0.092*** (0.026)	-0.052 (0.037)	-0.054 (0.049)	-0.360*** (0.093)
R2	0.888	0.904	0.806	0.917
N	249	249	249	249
Mean	3.71	4.89	6.65	15.48
% Mean/year	-2.48%	-1.06%	-0.81%	-2.32%
Tech FE	YES	YES	YES	YES
Year FE	YES	YES	YES	YES

Panel B:	<i>Normalized share<sub>i,τ,t</sub></i>			
	(1) Locations	(2) Industries	(3) Occupations	(4) All Firms
<i>Pioneer<sub>i,τ</sub></i>	2.393*** (0.528)	13.550*** (3.204)	10.746** (4.675)	142.036*** (35.866)
<i>Pioneer<sub>i,τ</sub> * years since emergence<sub>τ,t</sub></i>	-0.149*** (0.039)	-0.547** (0.224)	-0.367 (0.269)	-6.215** (2.990)
R2	0.076	0.137	0.033	0.026
N	266,467	26,883	204,041	38,990,627
$\beta(Pioneer_{i,\tau} * yrs\ since\ emg_{\tau,t}) / \beta(Pioneer_{i,\tau})$	-0.062*** (0.007)	-0.040*** (0.011)	-0.034*** (0.013)	-0.044*** (0.013)
Tech FE	YES	YES	YES	YES
Year FE	YES	YES	YES	YES
Cell FE	YES	YES	YES	YES

Notes: In Panel A, the table reports the results from a regression of the coefficient of variation calculated across  $Normalized\ share_{i,\tau,t} = \frac{share\ jobs\ exposed_{i,\tau,t}}{share\ jobs\ exposed_{\tau,t}}$  (where  $i$  is a location (Column 1), industry (Column 2), occupation (Column 3), or firm (Column 4)) for each technology  $\tau$  and time  $t$ . Location refers to a CBSA, industry is at the NAICS 4-digit level, and occupation is at the SOC 6-digit level. To keep the sample comparable across years, the coefficient of variation in Column (4) is calculated across 10,231 firms which have at least one job posting in each of the eleven years of Burning Glass; we call these ‘larger firms’. The specification excludes observations before the year of emergence a technology, and with fewer than 100 technology job postings which have an industry associated with them. Observations are weighted by the square root of total technology postings in a given year. The normalized share is capped at the 99<sup>th</sup> percentile of non-zero observations. Standard errors are clustered by technology. In Panel B, the results are from a regression of  $Normalized\ share_{i,\tau,t} = \frac{share\ jobs\ exposed_{i,\tau,t}}{share\ jobs\ exposed_{\tau,t}}$  on the Pioneer status of each cell (location/occupation/industry/firm), and its interaction with years since emergence. In Panel B column 4, we do not impose the balancing requirement, so that we have data on 300K unique firms in Burning Glass. The specification again excludes observations before the year of emergence a technology. Standard errors are clustered by cell. The expression  $\beta(Pioneer_{i,\tau} * years\ since\ emergence_{\tau,t}) / \beta(Pioneer_{i,\tau})$  is calculated by dividing the coefficient on  $Pioneer_{i,\tau} * Years\ since\ emergence_{\tau,t}$  by the coefficient on  $Pioneer_{i,\tau}$ . Standard errors for this expression are calculated using the delta method.



Table 11 - Robustness – Primary results for unsupervised approach

Panel A				
Dependent Variable:	Coefficient of Variation	Normalized Share	Share College Educated	Coefficient of Variation
Result:	Region Broadening (1)	Pioneer Persistence (2)	Skill Broadening (3)	Region Broadening by Skill (4)
<i>Years since emergence</i> <sub><math>\tau,t</math></sub>	--0.150*** (0.017)		--0.363*** (0.095)	--0.205*** (0.025)
<i>Pioneer</i> <sub><math>i,\tau</math></sub>		1.428*** (0.351)		
<i>Pioneer</i> <sub><math>i,\tau</math></sub> * <i>yrs since emg</i> <sub><math>\tau,t</math></sub>		--0.086*** (0.028)		
<i>Yrs since emg</i> <sub><math>skill,\tau,t</math></sub> * 1{ <i>skill</i> = low}				--0.094*** (0.029)
R2	0.841	0.023	0.873	0.691
N	2,339	2,135,310	2,339	4,663
Estimate (per year)	-2.56%	-6.23%	-0.68%	-1.03%
Mean	5.87	NA	53.07	9.17

Panel B				
	Coefficient of Variation			
	(1) Locations	(2) Industries	(3) Occupations	(4) Larger Firms
<i>Years since emergence</i> <sub><math>\tau,t</math></sub>	--0.150*** (0.017)	--0.032** (0.016)	--0.123*** (0.019)	--0.408*** (0.058)
R2	0.841	0.893	0.730	0.910
N	2,339	2,339	2,339	2,339
Mean(CV)	5.87	6.14	7.08	27.21
% Mean(CV)/year	-2.56%	-0.52%	-1.74%	-1.50%

Notes: This table reports our primary results replicated by treating each bigram as a separate technology (our ‘unsupervised approach’). In Panel A, we replicate our primary results in Table 4, Column 1, Table 5, Column 2, Table 6, Column 1, and Table 7, Column 3. In Panel B, we replicate results from Table 10, Panel A, again treating each of the 305 bigrams as a single technology. Observations are weighted by the square root of total technology job postings in a year. To keep the sample comparable across years, the coefficient of variation in Column (4) is calculated across 10,231 firms which have at least one job posting in each of the eleven years of Burning Glass; we call these ‘larger firms’. The normalized share is capped at the 99<sup>th</sup> percentile of non-zero observations. In Panel A, Columns 1, 3, and 4 control for technology and year fixed effects, while Column 2 controls for CBSA, technology and year fixed effects. In Panel B, all specifications control for technology and year fixed effects. Standard errors are clustered by technology in columns 1, 3, and 4 of Panel A, and in all columns of Panel B. Standard errors are clustered by CBSA in Column 2 in Panel A.

**Figure 1 – Sample job for Machine Learning/AI Technology**

**1.0.1 JobTitle**

Applied Research Scientist, Video Understanding

**1.0.2 JobText**

4.4 Facebook New York, NY Glassdoor Estimated Salary: 112k159k  
Applied Research Scientist, Video Understanding  
Facebook

Facebooks mission is to give people the power to build community and bring the world closer together. Through our family of apps and services, were building a different kind of company that connects billions of people around the world, gives them ways to share what matters most to them, and helps bring people closer together. Whether were creating new products or helping a small business expand its reach, people at Facebook are builders at heart. Our global teams are constantly iterating, solving problems, and working together to empower people around the world to build community and connect in meaningful ways. Together, we can help people build stronger communities were just getting started. Every day, massive amounts of video are uploaded into Facebooks services. In order to serve our communities better, it is critical that we can understand this content think about being able to answer questions like This person will like this video because.... or This person will find this video inappropriate because... Our goals broadly encompass content understanding, including the ability to produce video summaries, categorize content according to topic and purpose, identify audio events, find keyframes, and do keyword spotting. To achieve these goals, we are building a Video Understanding team in New York City, that will engage in a multidisciplinary effort combining speech recognition, natural language processing, and image processing. We view video as inherently multimodal content, and seek to develop methods that use all the information available. We are looking for researchers in machine learning and AI with strong software engineering skills, and a desire to build systems that will ship to billions of people. The Video Understanding Team is part of the Applied Machine Learning organization. The team carries out applied research in MLAI and designs, develops and deploys state of the art MLAI algorithms to the rest of Facebook. Our algorithms are used for ranking, improving content integrity, keeping communities safe, and power multiple product experiences across Facebook, Messenger, Instagram, WhatsApp and Oculus.

Responsibilities:

Develop highly scalable algorithms based on stateoftheart machine learning and neural network methodologies Conduct research to advance the stateoftheart, and publish work in relevant speech, NLP, and machine learning conferences and journals Apply expert coding skills to projects in partnership with other engineers across research, product, and infrastructure teams Adapt machine learning and neural network algorithms for training competitive, stateoftheart models which make the best use of modern parallel environments e.g. distributed clusters, GPU Minimum Qualifications:

MS degree in Computer Science or related quantitative field with 5 years of work experience, or Ph.D. degree in Computer Science or related quantitative field Knowledge of **machine learning, neural networks, and deep learning** Experience building systems based on machine learning and/or **deep learning** methods, especially in the areas of speech recognition, natural language processing, image processing, or other machineperception tasks Experience developing and debugging in CC and/or Python Preferred Qualifications:

Overview

Website [www.facebook.com](http://www.facebook.com)

Headquarters Menlo Park, CA, United States

Size 10000 employees

Founded 2004

Type Company Public FB

Industry Information Technology

**Figure 2– Sample job for Solar Technology**

**1.0.1 Job Title**

Solar Panel Installer

**1.0.2 JobDomain**

[www.glassdoor.co.in](http://www.glassdoor.co.in)

**1.0.3 JobDate**

20190315

**1.0.4 JobText**

3.4 Vivint Solar Corp Baltimore, MD

Solar Panel Installer

Vivint Solar Corp

**Job Description:** Right now, we are seeking a Solar Installer for our Rosedale MD office, who will be responsible for ensuring that our products are installed properly and ontime.

Must be have a valid drivers license Must be able to pass preemployment drug screen Must be able to pass criminal background

**Responsibilities:** Work with a team to install the racking system and solar panels on residential roofs Service the solar system as needed

**Required:** Working knowledge of solar installation, construction and/or roofing 1 to 2 years of relevant experience Ability to be comfortable being and working on roofs Valid drivers license Employees of Vivint Solar must submit to a criminal history check, motor vehicles check, drug screening, and obtain clearance from the state based upon the state requirements.

We do not accept resumes from headhunters, placement agencies, or other suppliers that have not signed a formal agreement with us.

Vivint Solar is a proud promoter of employment opportunities to our Military and Veterans. We, an equal opportunity employer, do not consider any protected traits e.g. race, creed, color, religion, gender, national origin, nonjobrelated disability, age, or any other protected trait when hiring under federal, state and local laws Company Description Vivint Solar is a leading fullservice residential solar provider in the United States.

With Vivint Solar, customers can power their homes with clean, renewable energy and typically achieve significant financial savings. Offering integrated residential solar solutions for the entire customer lifecycle, Vivint Solar designs, installs, monitors and services the solar energy systems for its customers. In addition to being able to purchase a solar energy system outright, customers may benefit from Vivint Solars affordable, flexible financing options or power purchase agreements. For more information, visit [www.vivintsolar.com](http://www.vivintsolar.com) or follow VivintSolar on Twitter.

Overview

Website [www.vivintsolar.com](http://www.vivintsolar.com)

Headquarters Lehi, UT, United States

Size 1001 to 5000 employees

Founded Unknown

Type Company Public VSLR

Industry Oil, Gas, Energy Utilities

Revenue 5 to 10 billion INR per year

Competitors Unknown

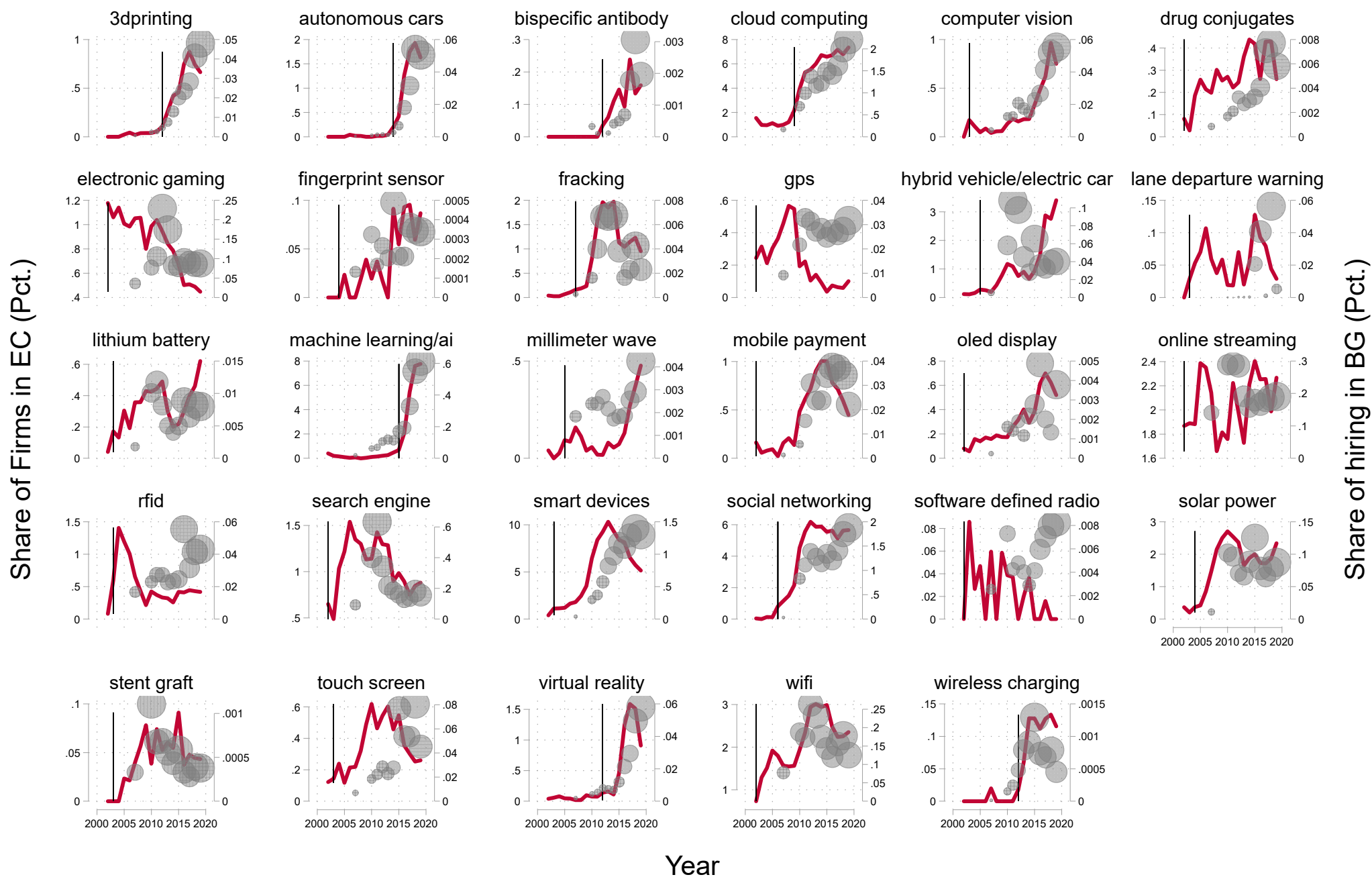
Vivint Solar Photos

Vivint Solar photo of: CEO, David Bywater, recognizing overachieving employees and thanking them for

**Notes:** Sample job posting which mentions Machine Learning/AI. The figure shows the standardized job title given by Burning Glass and the text of the job advertisement posted online on [glassdoor.com](http://glassdoor.com).

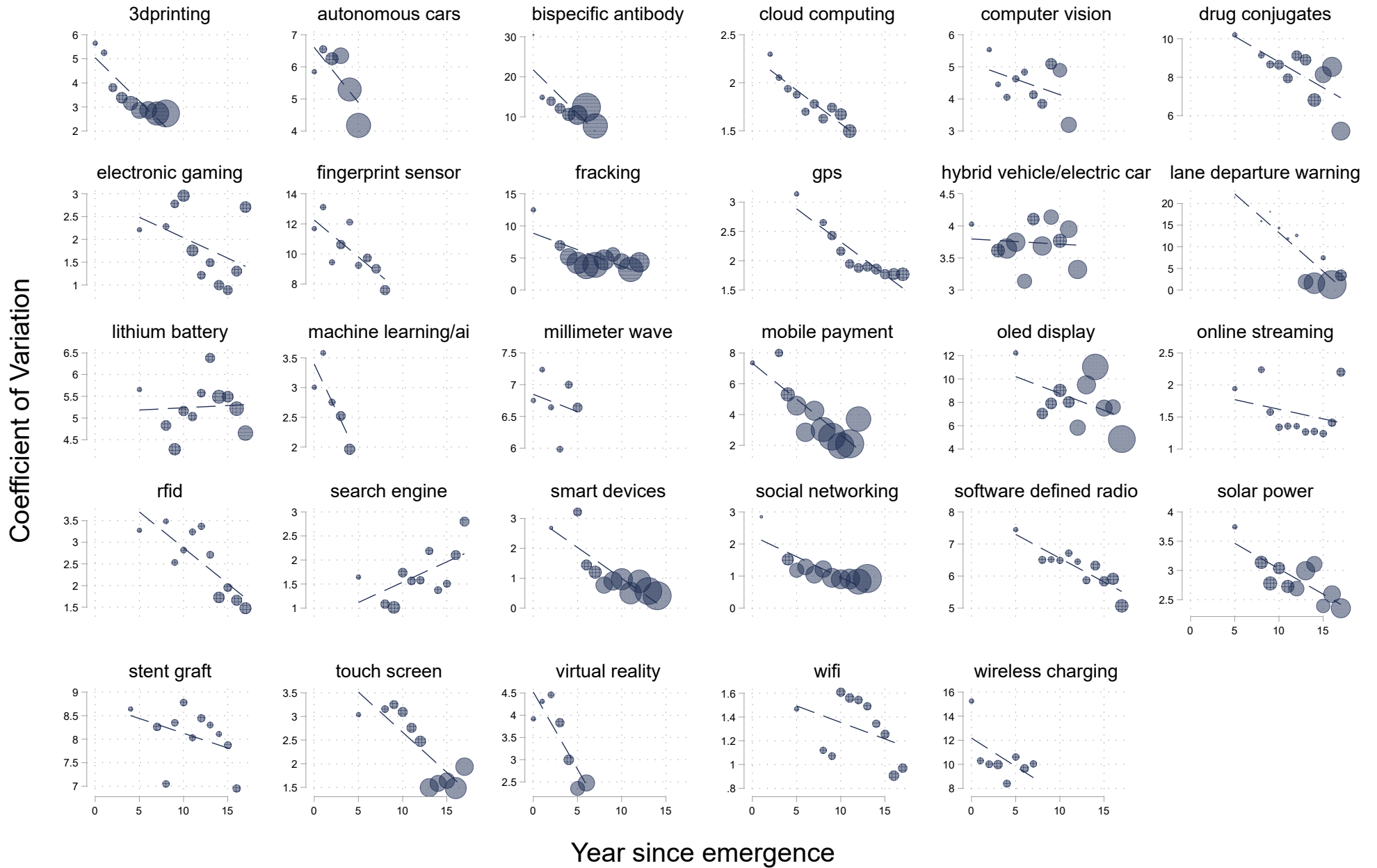
**Notes:** Sample job posting which mentions Solar technology. The figure shows the standardized job title given by Burning Glass and the text of the job advertisement posted online on [glassdoor.com](http://glassdoor.com).

**Figure 3 – Technology exposure in earnings calls and job postings, by year**



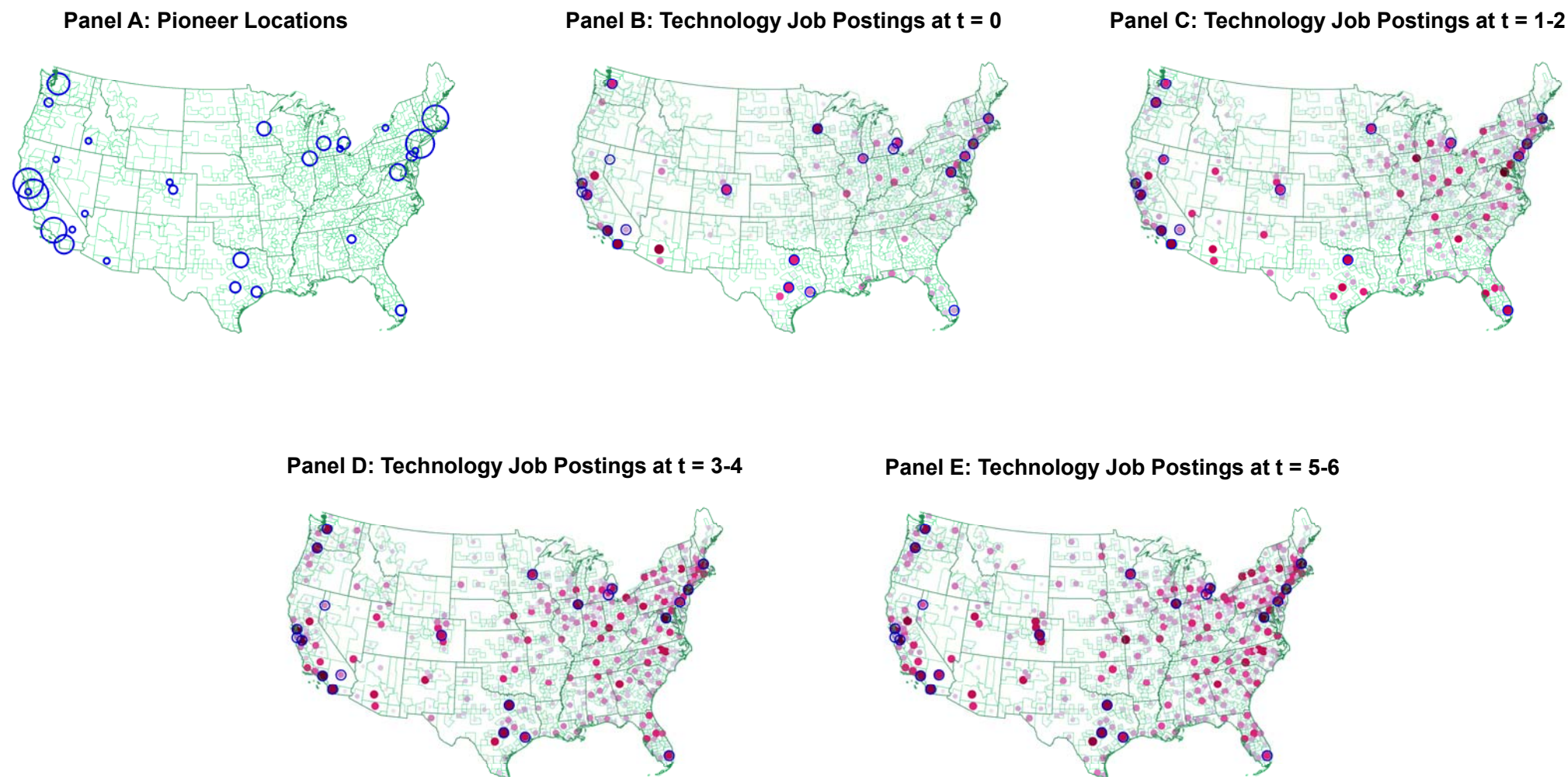
**Notes:** The figures plot (year by year) the percentage of firms (red line) that mention a given technology's keywords in earnings calls, and the percentage of job postings (grey circles) in Burning Glass that mention the same technology's keywords. The size of the grey circles denotes the level of postings for the technology x year observation. The vertical line highlights the year of emergence of the technology, which is defined as the year in which the earnings call time series (red line) attains at least 10% of the sample maximum. The overall correlation between these two time series is 81%.

**Figure 4 – Concentration of job-postings by geographic location, coefficient of variation by year since emergence**



**Notes:** The figure plots the coefficient of variation, measured using the normalized share of technology job postings for each of 29 technologies by year from 2007 to 2019, where  $Normalized\ share_{i,\tau,t} = \frac{share\ jobs\ exposed_{i,\tau,t}}{share\ jobs\ exposed_{\tau,t}}$ , for CBSA  $i$ , years since emergence  $t$ , and technology  $\tau$ . Only observations at the year of emergence and after are included.

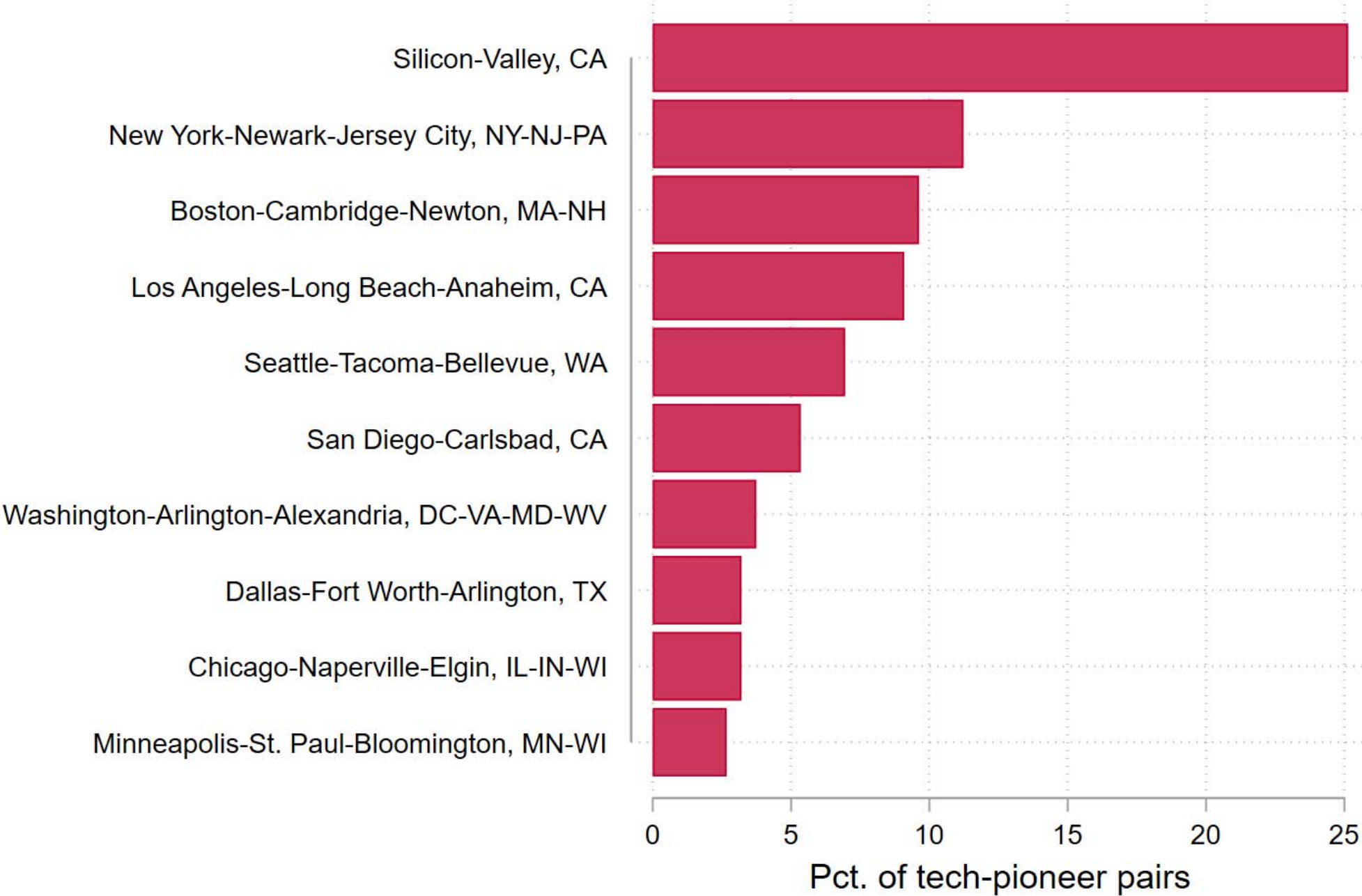
**Figure 5 – Geographic diffusion of job postings exposed to disruptive technologies, by year since emergence**



**Notes:** In Panel A, we show as circles the Core-Based Statistical Areas (CBSAs) that are pioneer locations for at least one technology. The size of the circles is proportional to the share of technologies for which the CBSA is a pioneer location. Panels B through E continue to mark pioneer locations with hollow blue circles, but now also add the location of technology job postings in the start year of the technology (the average  $Normalized\ share_{i,\tau,t}$  across technologies at  $t=0$ ), where darker dots correspond to a higher normalized share of job postings. We plot these pictures only for the 13 out of our 29 technologies that have a year of emergence after 2007, so that the panel remains balanced throughout.

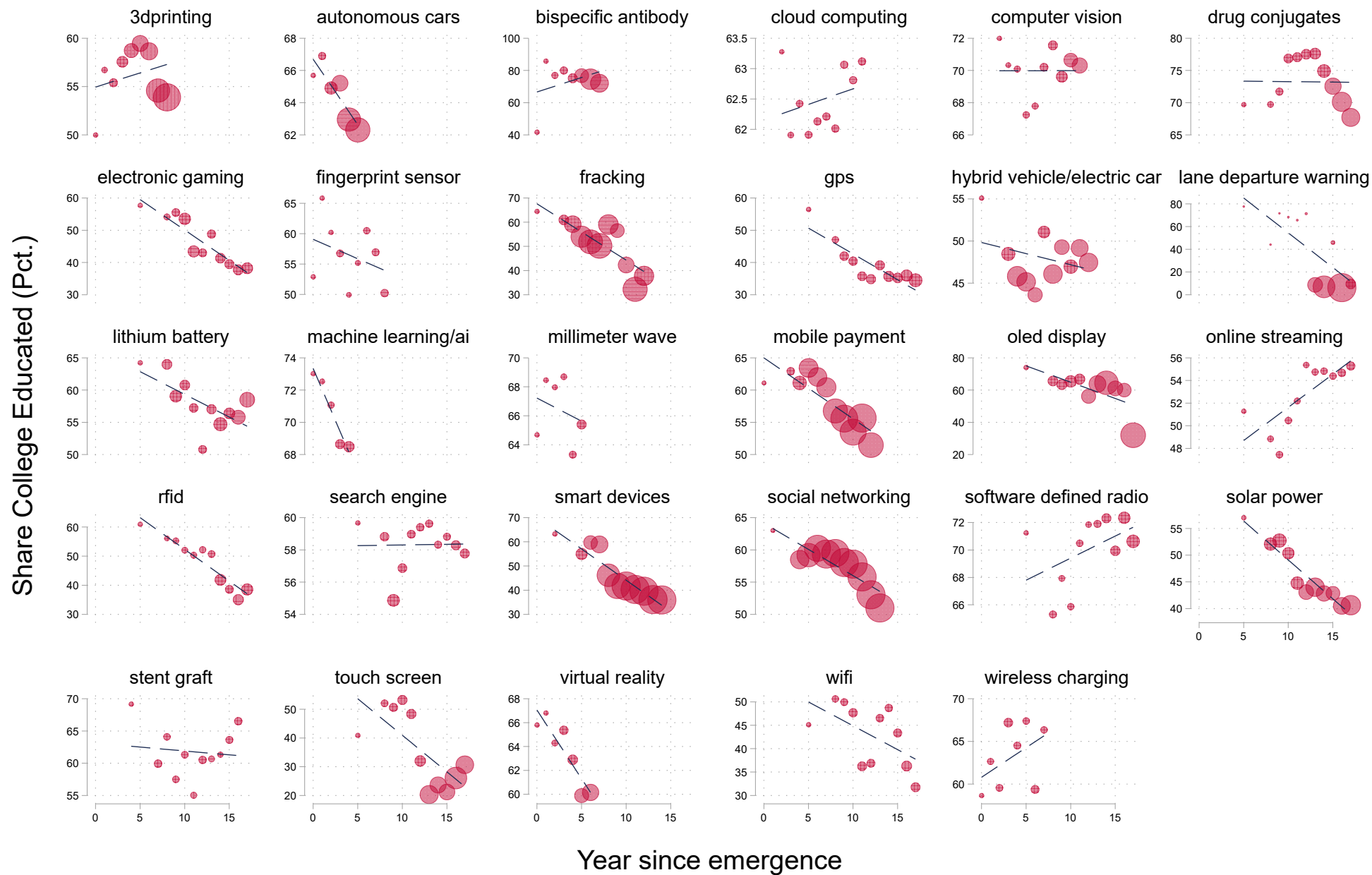


**Figure 6 – Percentage of technology-pioneer location pairs, by region**



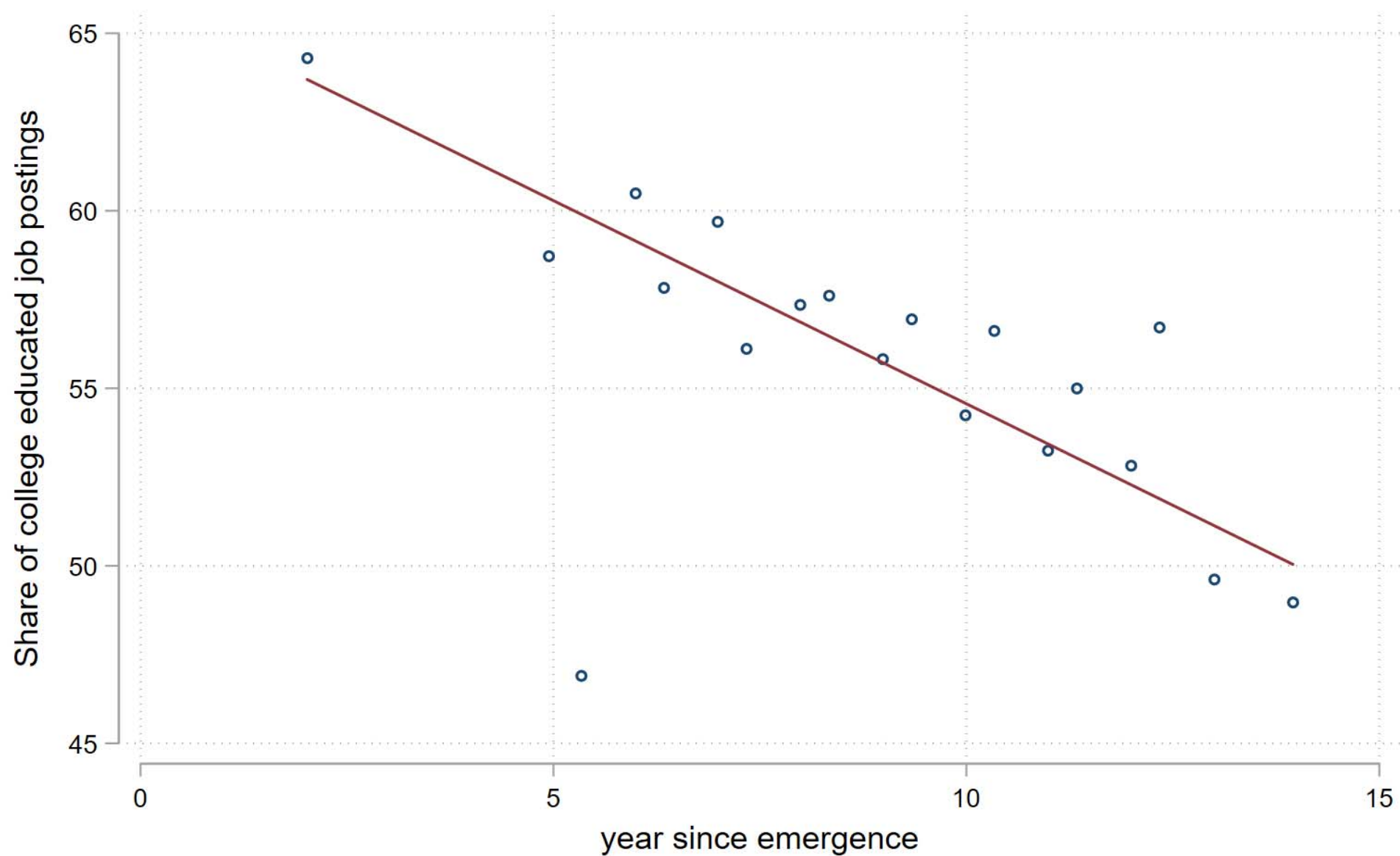
**Notes:** The figure plots the percentage of technology-pioneer location pairs by region for the top 10 regions. The total number of technology-pioneer locations is calculated as the number of technologies for which CBSAs in the region are classified as pioneer locations. We combine San Jose-Sunnyvale-Santa Clara, CA and San Francisco-Oakland-Hayward, CA, and jointly label the region as Silicon Valley. Jointly, all California CBSAs account for 40.2% of all technology-pioneer location pairs. Major cities in the Northeast corridor, New York-Newark-Jersey City, NY-NJ-PA, Boston-Cambridge-Newton, MA-NH, and Washington-Arlington-Alexandria, DC-VA-MD-WV, jointly account for 24.58% of all total technology-pioneer location pairs.

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



**Notes:** The figure plots the approximate share of technology job postings that require a college education (red circles, where the size of the circle represents the total number of technology job postings) over the years since the emergence of the technology. The approximate share of technology job postings that require a college education is calculated using  $Skill_t^{\tau} = \frac{\sum_o N_{o,t}^{\tau} \chi_o}{\sum_o N_{o,t}^{\tau}}$ , where  $\chi_o$  is the college share in an occupation in the 2015 American Community Survey and  $N_{o,t}^{\tau}$  is the number of technology job postings in technology  $\tau$ . Only observations at the year of emergence and after are included.

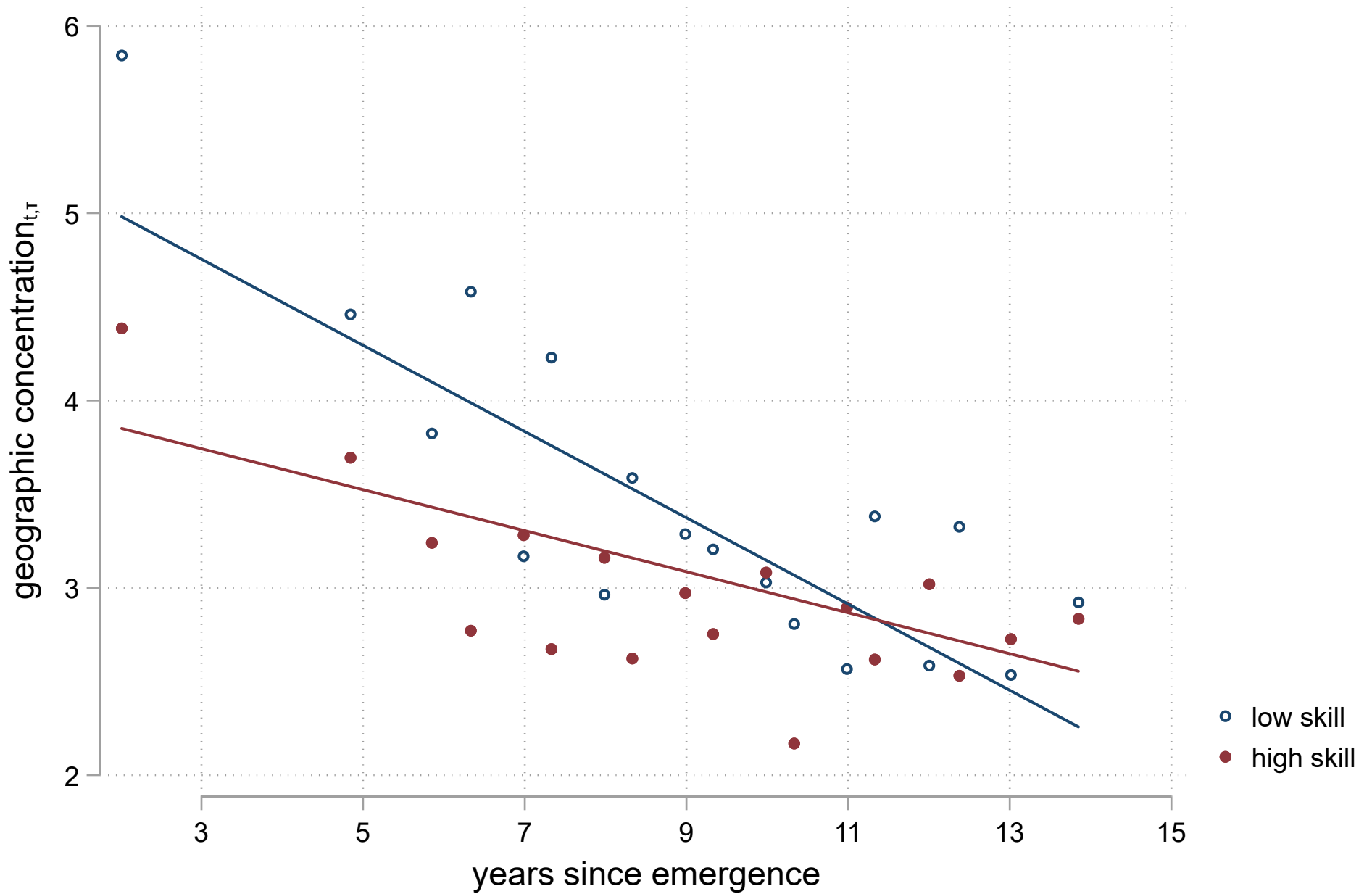
**Figure 8 - Share of job postings requiring college education pooled across technologies, by year since emergence**



**Notes:** The figure shows a bin scatter of the approximate share of technology job postings requiring a college education against the years since emergence of technology. The sample pools across technologies where each observation in the sample denotes a technology x year observation. We weight observations by the square root of job postings in that technology x year pair. Only observations at the year of emergence and after are included. The fitted line controls for technology fixed effects.

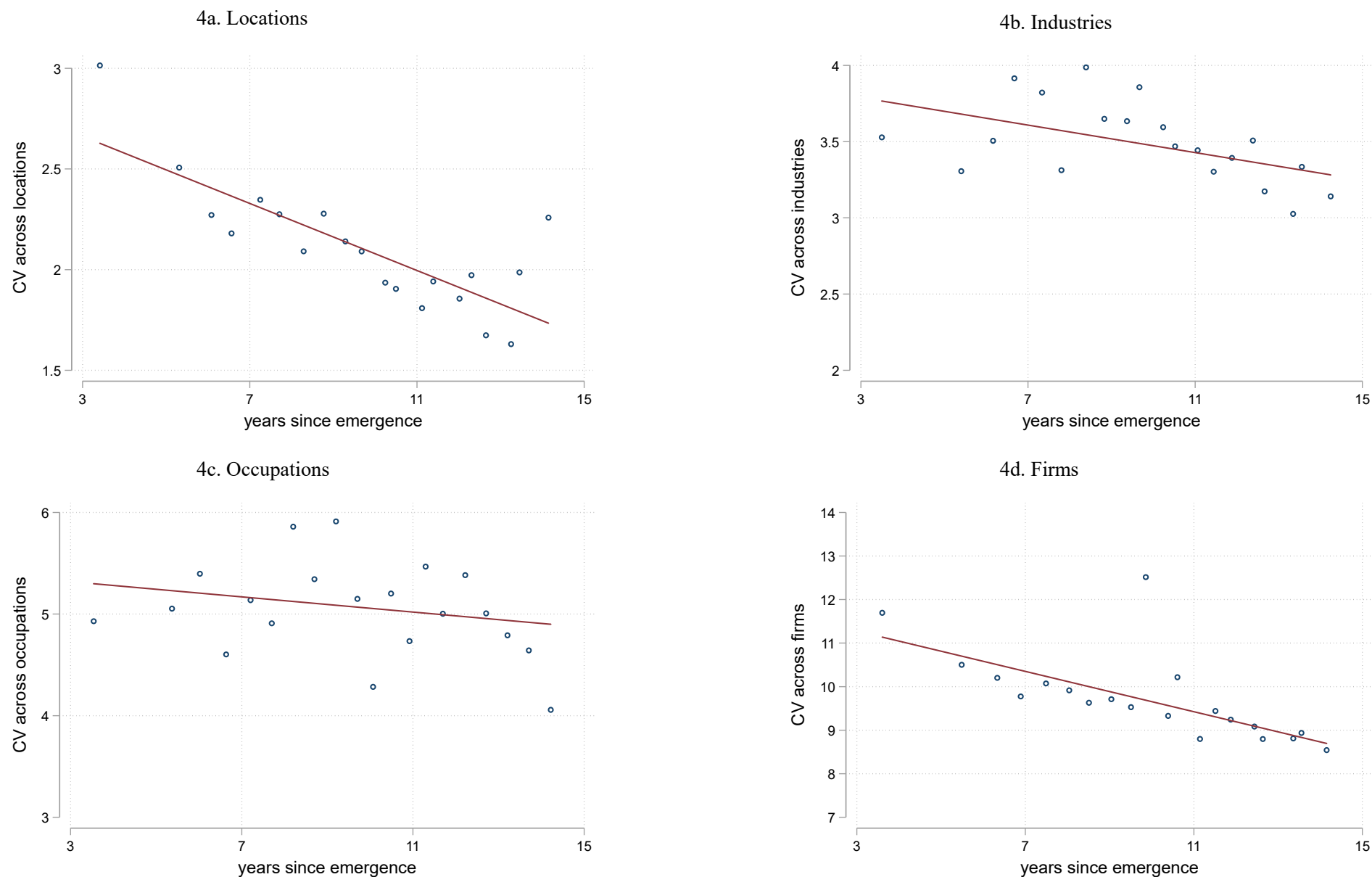


**Figure 9 – Geographic concentration by skill, coefficient of variation by year since emergence**



Notes: This figure shows a binned scatter plot of the coefficient of variation of the *Normalized share* of technology job postings against the years since the emergence of the technology, separately for high-skill (red) and low-skill (blue) occupations. For a given technology, year, and skill level, we calculate the coefficient of variation of *Normalized share*<sub>*i*,*τ*,*t*,*s*</sub> =  $\frac{\text{share jobs exposed}_{i,\tau,t,s}}{\text{share jobs exposed}_{\tau,t,s}}$ , where *i* is a CBSA, *t* denotes years since emergence of the technology, and *s* is the skill level. The fitted lines control for technology fixed effects and weigh observations by the square root of the total number of postings in the technology-year pair. High-skill occupations are those with a share of college-educated people > 60% in the 2015 ACS and low-skill occupations are those with a share of college-educated people < 30%. *Normalized share* is capped at the 99<sup>th</sup> percentile of non-zero observations.

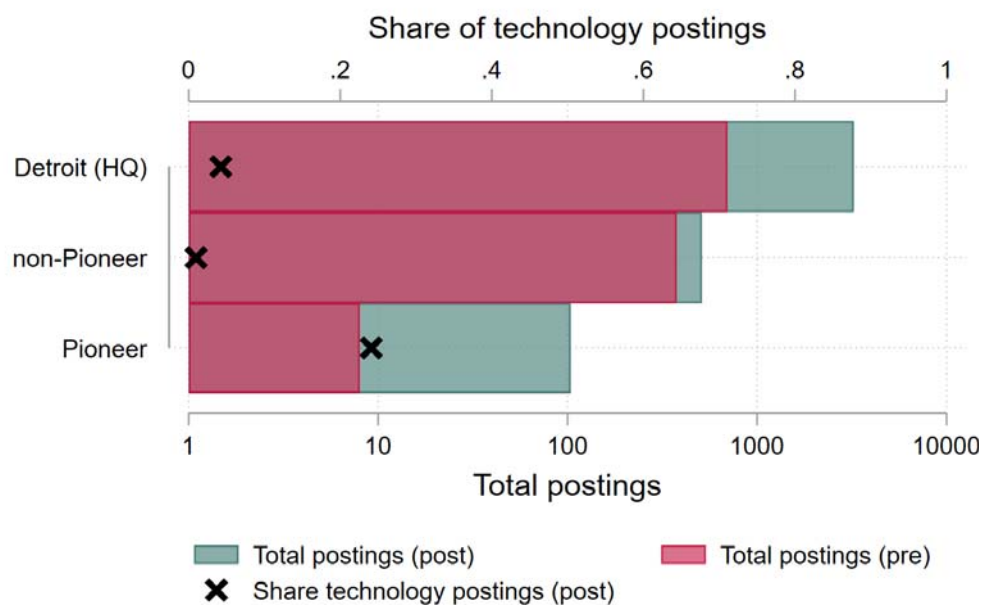
**Figure 10 – Geographic concentration of job postings across locations, industries, occupations, and firms**



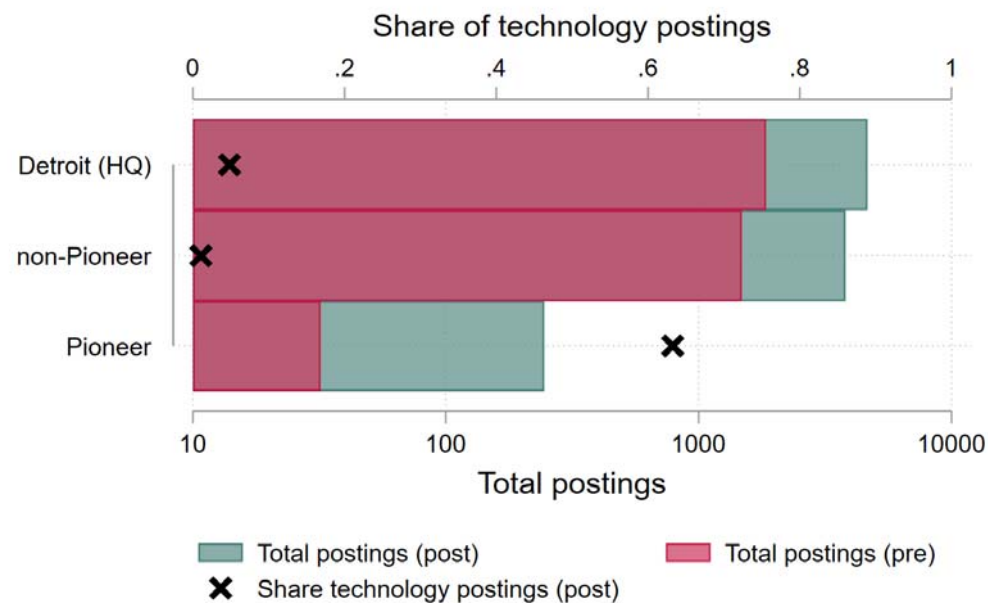
Notes: The figures are binned scatter plots of the coefficient of variation over the *years since emergence* for our panel of 29 technologies. The coefficient of variation is calculated across  $Normalized\ share_{i,\tau,t} = \frac{share\ jobs\ exposed_{i,\tau,t}}{share\ jobs\ exposed_{\tau,t}}$  (where  $i$  is a location (4a), industry (4b), occupation (4c), or firm (4d)) for each technology  $\tau$  and time  $t$ . Location refers to a CBSA, industry is at the NAICS 4-digit level, and occupation is at the SOC 6-digit level. The coefficient of variation in Figure (4d) is calculated across 10,231 firms that have at least one job posting in each of the 11 years of our Burning Glass data. The results only include observations at the time of and after the emergence year of a technology, and observations with more than 100 technology job postings that have an industry associated with them. The binscatter is weighted by the square root of total technology postings in a year. *Normalized share* is capped at the 99<sup>th</sup> percentile of non-zero observations.

**Figure 11 – Rehoming of firms to Pioneer locations**

### Ford Motor Company



### General Motors Corp.



**Notes:** The figure decomposes job postings of Ford Motor Company and General Motors Corp between their headquarters (Detroit), pioneer CBSAs for the Autonomous Cars technology, and other CBSAs (non-Pioneer). The bars decompose all job postings while the x's decompose technology job postings. The level of red bars indicates all postings before the year of emergence of the Autonomous Cars technology and the level of green bars indicate all postings after the year of emergence. The level of x's indicate share of technology job postings in overall job postings post year since emergence. The decompositions are calculated separately for Ford Motor Company (on the left) and General Motor Corporation (on the right). There are four pioneer locations for the Autonomous Car technology: San Jose-Sunnyvale-Santa Clara (CA), San Francisco-Oakland-Hayward (CA), Boston-Cambridge-Newton (MA-NH), Detroit-Warren-Dearborn (MI). Headquarters for Ford and GM are located in Detroit-Warren-Dearborn (MI), which is also a pioneer location for Autonomous Cars technology.

## Appendix Tables and Figures

Appendix Table 1 – Alternative cutoffs defining disruptive technologies

(1)	(2)	(3)	(4)	(5)
$\frac{EC_{2002}}{EC_{max}} \leq 0.01$	$\frac{EC_{2002}}{EC_{max}} \leq 0.05$	$\frac{EC_{2002}}{EC_{max}} \leq 0.1$	$\frac{EC_{2002}}{EC_{max}} \leq 0.15$	$\frac{EC_{2002}}{EC_{max}} \leq 0.2$
total bigrams: 235	total bigrams: 243	total bigrams: 305	total bigrams: 401	total bigrams: 500
machine learning	machine learning	mobile devices	mobile devices	mobile devices
cloud computing	cloud computing	machine learning	machine learning	financial instruments
cloud services	cloud services	cloud computing	cloud computing	machine learning
public cloud	public cloud	cloud services	solid organic	cloud computing
social networking	social networking	quality metrics	cloud services	elevated levels
smart grid	smart grid	flow profile	data usage	solid organic
cloud service	cloud service	smart phones	quality metrics	cloud services
cloud infrastructure	cloud infrastructure	mobile platform	flow profile	data sets
carbon footprint	carbon footprint	public cloud	smart phones	data usage
virtual reality	virtual reality	social networking	video content	quality metrics
autonomous driving	social networks	smart grid	mobile platform	flow profile
augmented reality	autonomous driving	cloud service	public cloud	smart phones
cloud environment	augmented reality	connected devices	social networking	video content
autonomous vehicles	cloud environment	cloud infrastructure	smart phone	mobile platform
cloud based	autonomous vehicles	carbon footprint	lifecycle management	public cloud
hydraulic fracturing	global warming	nand flash	smart grid	social networking
wifi network	cloud based	virtual reality	cloud service	smart phone
results page	hydraulic fracturing	digital channel	connected devices	fleet management
additive manufacturing	optimization process	delivery network	mobile platforms	lifecycle management
relevant content	software defined	social networks	cloud infrastructure	smart grid

Notes: For each alternative cutoff (listed in Row 1), the table reports the number of bigrams in Row 2 and the set of top 20 most frequent bigrams in earnings calls from Row 3 onward. Each cutoff is calculated as the ratio of the share of earnings calls mentioning the bigram in 2002 and the maximum between 2002 and 2019 of the share of earnings calls mentioning the bigram (denoted by  $\frac{EC_{2002}}{EC_{max}}$ ). In the successive columns, we keep bigrams that increase in their frequency of use by factors 100, 20, 10 (our baseline specification), 6.66, and 5 during our sample period, respectively.

Appendix Table 2 – Alternative cutoffs for minimum number of mentions in earnings calls

(1)	(2)	(3)	(4)	(5)
$EC_{total} \geq 80$	$EC_{total} \geq 90$	$EC_{total} \geq 100$	$EC_{total} \geq 110$	$EC_{total} \geq 120$
Total bigrams: 406	Total bigrams: 347	Total bigrams: 305	Total bigrams: 279	Total bigrams: 240
Bigrams w/ $89 \geq EC_{total} \geq 80$	Bigrams w/ $99 \geq EC_{total} \geq 90$	Bigrams w/ $109 \geq EC_{total} \geq 100$	Bigrams w/ $119 \geq EC_{total} \geq 110$	Bigrams w/ $129 \geq EC_{total} \geq 120$
software element	cell activation	digital color	cloud server	response based
object storage	retinal vein	phase production	injection molded	virtual environment
content creator	mobile telephones	nand memory	sensor fusion	keyless entry
sense multiple	pharmacokinetic properties	impact analysis	illustration purposes	primary fuel
recycled fibers	docking station	diesel fuels	mobility management	optimization system
search algorithms	optic fiber	specific network	activation process	advertising messages
communication modules	tissue sarcoma	system failures	frequency identification	touch sensor
controlled process	system enabling	laser scanner	neural network	grinding media
user activities	receptor agonists	launch system	charging station	lateral section
bile acid	vehicle components	cancer stem	conversion efficiencies	fingerprint sensor
expected data	coated product	resource pool	virtual currency	load management
vehicle tires	distillation unit	lumbar spine	direct fuel	search result
comparison results	patient monitor	whey protein	purchase history	tumor samples
cell survival	computed tomography	treatment cycle	fluid volumes	gaming environment
allergic conjunctivitis	computing resources	reverse circulation	lupus nephritis	brake pads
networking site	bispecific antibodies	service layer	electronic gaming	esophageal cancer
cancer melanoma	network content	resource intensive	centralized data	target volume
flow measurement	sustained delivery	multiple mobile	driving behavior	imaging modalities
network communication	ethylene propylene	inkjet printing	search tool	data connections
physical channels	target regions	database technology	thermal capacity	primary storage

Notes: For varying earnings calls cut-offs (mentioned in Row 1), the table reports the number of bigrams in Row 2 and the most frequently mentioned “marginal” bigrams from Row 3 onward. Marginal bigrams in Column 1 are ones which are mentioned in 80-89 earnings calls, and hence get dropped at the margin when switching to a more stringent requirement of at least 90 mentions (in Column 2). For example, “software element” is one of these marginal bigrams mentioned in more than 80 earnings calls but less than 90. Column 3 shows our baseline specification corresponding to Table 1.

Appendix Table 3 – Technology descriptions and contemporaneous events around emergence year

Technology	Description	Year of emergence	Contemporaneous Event
Smart devices	A smart device is an electronic device, generally connected to other devices or networks via different wireless protocols such as Bluetooth, Zigbee, NFC, Wi-Fi, LiFi, 5G, etc., that can operate to some extent interactively and autonomously.	2005	Apple announces first iPad. – Apple (2005)
Cloud computing	Cloud computing is the on-demand availability of computer system resources, especially data storage and computing power, without direct active management by the user.	2008	Microsoft and Google announced their cloud platforms. – Google and Microsoft blogs (2008)
Social networking	The use of dedicated websites and applications to interact with other users, or to find people with similar interests to oneself.	2006	Mark Zuckerberg leaves Harvard. – Harvard Crimson (2005) Facebook receives \$25 mill venture funding, and valued at half a billion. – Market Watch (2006)
Machine learning/AI	Machine learning focuses on the development of computer programs that can access data and use it to learn for themselves.	2015	Tesla's Elon Musk and venture capitalist Peter Thiel dedicated \$1 billion to found Open AI, a non-profit for artificial intelligence research. – USA Today (2015)
Solar power	Solar power is the conversion of energy from sunlight into electricity, either directly using photovoltaics (PV), indirectly using concentrated solar power, or a combination.	2002	Gov. Arnold Schwarzenegger announces plans for solar power subsidies. – Sacramento Bee (2005)
Autonomous cars	A self-driving car, also known as an autonomous vehicle (AV), connected and autonomous vehicle (CAV), full self-driving car or driverless car, or robo-car or robotic car, (automated vehicles and fully automated vehicles in the European Union) is a vehicle that is capable of sensing its environment and moving safely with little or no human input.	2014	Google unveiled its first "fully functional" prototype for a self-driving car Monday and plans to test it on Bay Area public roads in the new year. – Mercury News, The (2014)
Virtual reality	Virtual reality (VR) refers to a computer-generated simulation in which a person can interact within an artificial three-dimensional environment using electronic devices, such as special goggles with a screen or gloves fitted with sensors.	2013	Oculus raises \$16 million in venture funding for virtual reality headset. – The Verge (2013)
Search engine	A search engine is a software system that is designed to carry out web searches (Internet searches), which means to search the World Wide Web in a systematic way for particular information specified in a textual web search query.	2002	Sausalito start-up Groxis released a new search tool that categorizes search results in a more visually friendly way. - Mercury News, The (2003)
Hybrid vehicle/Electric car	Any land-based automobile which uses electricity as one of the power sources.	2007	Toyota announces its plans for a plug-in hybrid car. – New York Times (2008). The Obama Administration lends Tesla Motors \$465 million to build an electric sedan and the battery packs needed to propel it. – Wired (2009)
Wireless charging	Inductive charging (also known as wireless charging or cordless charging) is a type of wireless power transfer. It uses	2012	General Motors invest \$5 million in wireless charging start-up Powermat. – Reuters (2012)

	electromagnetic induction to provide electricity to portable devices.		
Touch screen	The touchscreen enables the user to interact directly with what is displayed, rather than using a mouse, touchpad, or other such devices (other than a stylus, which is optional for most modern touchscreens).	2003	Santa Clara county uses touch machines for voting. San Jose Mercury News (2003)
Drug conjugates	Antibody-drug conjugates or ADCs are a class of biopharmaceutical drugs designed as a targeted therapy for treating cancer.	2002	Seattle Genetics signed a licensing deal granting MedImmune rights to use its antibody-drug-linking technology in research against a single biological marker of cancer. – Seattle Times, The (2005)
Fracking	Hydraulic fracturing, also called fracking, fracing, hydrofracking, fraccing, frac'ing, and hydrofracturing, is a well stimulation technique involving the fracturing of bedrock formations by a pressurized liquid.	2007	Congress signs fracking as an exemption from the Safe Drinking Water Act. - Denver Post, The (CO) (2003)
Software defined radio	Software-defined radio (SDR) is a radio communication system where components that have been traditionally implemented in hardware (e.g. mixers, filters, amplifiers, modulators/demodulators, detectors, etc.) are instead implemented by means of software on a personal computer or embedded system.	2002	Boeing was awarded a \$220 million subcontract to Northrop Grumman's Radio Systems business in San Diego to expand development of the communications, navigation and identification system specializing in software-defined radios for the Army's Comanche helicopter. San Diego Union-Tribune, The (2003)
Wi-Fi	Wi-Fi is a family of wireless network protocols, based on the IEEE 802.11 family of standards, which are commonly used for local area networking of devices and Internet access.	2002	San Francisco officials invited responses from 17 companies - including Google - that are interested in bringing affordable wireless Internet connections to the entire city. – Mercury News, The (2005)
3D printing	3D printing, or additive manufacturing, is the construction of a three-dimensional object from a CAD model or a digital 3D model.	2011	Federal government released plans to spend \$45 million to help fund a planned additive manufacturing institute. - USA Today (2012)
Millimeter wave	Extremely high frequency (EHF) or Millimeter Wave is the International Telecommunication Union (ITU) designation for the band of radio frequencies in the electromagnetic spectrum from 30 to 300 gigahertz (GHz).	2014	Facebook develops millimeter-wave networks for Internet.org. – The Verge (2016)
GPS	The Global Positioning System (GPS), originally NAVSTAR GPS,[1] is a satellite-based radionavigation system owned by the United States government and operated by the United States Space Force.[2]	2002	The Clinton administration removes “Selective Availability” of civilian GPS in order to make it more useful worldwide. – GPS.gov (2000)
Lithium-ion battery	A lithium-ion battery or Li-ion battery is a type of rechargeable battery.	2002	Sion Power Corp. starts production of a new lithium-sulfur battery that can last twice as long as the previous model commonly used in laptops, cell phones and digital cameras. - Arizona Daily Star, The (2004)

OLED display	An organic light-emitting diode (OLED or organic LED), also known as organic electroluminescent (organic EL) diode, is a light-emitting diode (LED) in which the emissive electroluminescent layer is a film of organic compound that emits light in response to an electric current.	2002	Kodak announced the first consumer product to include a full-color, active-matrix organic light-emitting diode (OLED) display on the Kodak EasyShare LS633 digital camera. - Mercury News, The (2003)
Stent graft	In medicine, a stent is a metal or plastic tube inserted into the lumen of an anatomic vessel or duct to keep the passageway open, and stenting is the placement of a stent.	2003	A stent graft system designed to correct life-threatening thoracic aortic aneurysms is fast track approved by the Food and Drug Administration. - Houston Chronicle (2003)
RFID	Radio-frequency identification (RFID) uses electromagnetic fields to automatically identify and track tags attached to objects. An RFID tag consists of a tiny radio transponder; a radio receiver and transmitter.	2002	Wal-Mart Stores ordered its 100 top suppliers to begin using RFID tags on shipments beginning in January 2005. - Mercury News, The (2003)
Electronic gaming	An electronic game is a game that employs electronics to create an interactive system with which a player can play.	2002	Sony launches PlayStation 2 capable of playing video games from DVDs. Gamespy.com (1999) Microsoft launches Xbox, first mainstream device with online capabilities. Xbox.com (2000)
Computer vision	Computer vision is an interdisciplinary scientific field that deals with how computers can gain high-level understanding from digital images or videos.	2003	The state of Illinois processes 10 million driver's license images using facial recognition. Chicago Sun-Times (2002)
Lane departure warning	In road-transport terminology, a lane departure warning system (LDWS) is a mechanism designed to warn the driver when the vehicle begins to move out of its lane (unless a turn signal is on in that direction) on freeways and arterial roads.	2002	Iteris and DaimlerChrysler develop a first Lane Departure Warning System. The device is mounted on a truck's windshield. It houses a tiny camera, computer and software that tracks the difference between the road and visible lane markings. Seattle Times, The (2003)
Bispecific monoclonal antibody	A bispecific monoclonal antibody (BsMAb, BsAb) is an artificial protein that can simultaneously bind to two different types of antigen. BsMabs can be manufactured in several structural formats, and current applications have been explored for cancer immunotherapy and drug delivery.	2012	Novartis Pays Genmab \$2M to research into Bispecific Antibody Technology - Genetic Engineering and Biotechnology News (June 2012)
Fingerprint sensor	A fingerprint sensor is an electronic device used to capture a digital image of the fingerprint pattern.	2011	Apple buys fingerprint sensor firm AuthenTec for \$356 million. - ZDNet (July 2012)
Mobile payment	Mobile payment (also referred to as mobile money, mobile money transfer, and mobile wallet) generally refer to payment services operated under financial regulation and performed from or via a mobile device.	2007	Bank of America, Citibank, Wachovia, Washington Mutual, Wells Fargo, and ING Direct announce mobile banking services, including mobile payment services. - CNBC (June 2007)
Online streaming	Streaming media is multimedia that is constantly received by and presented to an end-user while being delivered by a provider.	2002	Apple invested \$12.5 million in Akamai, a content delivery company, with the aim to develop video streaming services for QuickTime TV. Akamai.com (1999)

Notes: The table above lists our 29 technologies (in Column 1), definitions for each technology taken from Wikipedia (Column 2), the emergence year based on earnings calls (Column 3), and a suggested contemporaneous event around the year of emergence (Column 4).



Appendix Table 4 - Example of keyword human audit – “Autonomous car” keywords

Bigrams	True Positive Rate	Comments	Status
autonomous vehicle*	100%		Keep
autonomous vehicles*	100%		Keep
autonomous driving*	100%		Keep
self-driving car	90%		Keep
automated car	0%	-automated car washes.	Drop
robotic car	0%	- robotic car wash, - "robotics, car," - Shelley the robotic car from a video	Drop
robot car	100%		Keep
driverless car	90%		Keep
driverless truck	100%		Keep
autonomous car	100%		Keep
driver assistance	0%	- [role of a] senior living team members who is performing in a <b>driver assistance</b> role, spotter, or resident care.	Drop
automated driving	100%		Keep
autonomous cars	100%		Keep

Notes: The table presents the results for a human audit of keywords for the “autonomous car” technology. For the audit, we go through each of the shortlisted bigrams (in Column 1) and randomly sample 10 job postings. Column 2 presents the true positive rate and Column 3 shows comments from auditor. Column 4 shows whether we keep or drop the bigram for the final list. \* Bigrams are ones which we originally obtained from the intersection of patent corpus with earnings calls. The remaining bigrams were short-listed during the auditing process.

Appendix Table 5 – Technology keywords

Technology	Keywords
3D printing	3d printer; 3d printing; additive manufacturing; 3d printed
Autonomous cars	Self-driving car; robot car; autonomous vehicles; autonomous car; autonomous cars; automated driving; driverless car; autonomous driving; autonomous vehicle; driverless truck
Bispecific monoclonal antibody	bispecific monoclonal; the bispecific; bispecific antibody
Cloud computing	paas; cloud infrastructure; distributed cloud; cloud provider; cloud offerings; cloud service; cloud applications; community cloud; private cloud; public cloud; cloud deployments; cloud environments; cloud management; cloud services; cloud security; enterprise class; iaas; hybrid cloud; cloud platform; cloud providers; cloud hosting; personal cloud; enterprise network; cloud computing; cloud based; saas; cloud storage; enterprise applications; cloud solution; enterprise cloud; cloud solutions; cloud deployment
Computer vision	pose estimation; motion estimation; visual servoing; facial recognition; gesture recognition; computer vision; image recognition; sensor fusion; object recognition
Drug conjugates	kinase inhibitor; drug conjugate; antibody drug; drug conjugates
Electronic gaming	social game; video games; social games; video game; game content; electronic gaming; gaming products
Millimeter wave	millimeter wave
Fingerprint sensor	fingerprint sensor; fingerprint scanner
Fracking	fracking; fraccing; hydrofracking; hydrofracturing; hydraulic fracturing
GPS	gps systems; global positioning; navigation devices
Hybrid vehicle/Electric car	hybrid vehicle; electric vehicle; electric motorcycle; vehicle charging; hybrid electric; plugin hybrids; electric buses; electrical vehicles; electric car; electric vehicles
Lane departure warning	lane departure; departure warning
Lithium battery	ion battery; lithium ion battery; lithium ion batteries; lithium batteries; ion batteries; lithium polymer; lithium ion; lithium battery
Machine Learning/AI	neural network; deep learning; language processing; machine learning; machine intelligence; natural language; artificial intelligence; AI technology; supervised learning; learning algorithms; unsupervised learning; reinforcement learning; AI machine
Mobile payment	mobile transfer; mobile commerce; mobile payment; mobile wallet; mobile money
OLED display	oled
Online streaming	streaming content; music streaming; interactive tv; live stream; digital video; video conferencing; online streaming; online video; mobile video; streaming services; streaming media; live video; video ondemand; live streaming; video ad; internet radio; video streaming; streaming video
RFID	frequency identification; keyless entry; rfid tags; rfid
Search Engine	search engine; search engines
Smart devices	mobile devices; tablet computers; wearable devices; tablet pcs; smartphone tablet; android phones; media devices; smart phones; smart devices; smart tvs; smart speaker; smart watch; smart car; smart phone; iphone ipad; portable media; smart tablets; connected devices; smartphones tablets; android smartphones; phones tablets; android devices; smart refrigerator; smartcar; smartphone; smart tv; smart band
Social networking	user generated; user generated content; social platforms; networking sites; social channels; social media; social networking; social networks; social network
Software defined radio	defined radio
Solar Power	solar wafer; rooftop solar; solar modules; solar cells; crystalline silicon; silicon solar; solar panel; solar power; solar wafers; solar energy; solar applications; solar module; solar cell; solar pv; solar grade; solar panels; photovoltaic; solar thermal
Stent graft	stent graft
Touch screen	touch controller; touch panel; capacitive touch; touchscreen; touch screens; touch sensor
Virtual reality	virtual reality; augmented reality; mixed reality; extended reality
Wi-Fi	wifi hotspots; wifi network; wifi; broadband connectivity; wireless networks
Wireless charging	wireless charging; inductive charging

Notes: This table shows, for each of our 29 technologies (in Column 1), the full set of 221 keywords used to associate earnings calls, patents and job postings with the technology.

Appendix Table 6 - Technology Excerpts from earnings calls

Company	EC month	Excerpt
Ambarella Inc	4/2018	results that are many times higher in terms of processing performance per watt In March we successfully demonstrated to customer and investors our fully  AUTONOMOUS VEHICLE or embedded vehicle autonomy on Silicon Valley Road EVA navigated various traffic scenarios presented by Silicon Valleys challenging urban environment The fully autonomous
General Motors Co	7/2017	safely deploy our selfdriving electric vehicles in commercial ridesharing networks Last month GM became the first company to use mass production methods to build  AUTONOMOUS VEHICLES  growing our test fleet to We plan to deploy these vehicles in the challenging driving environment of San Francisco as well as Scottsdale Arizona
Agenus Inc	10/2019	differentiated proof of mechanism of our potentially first or bestinclass agents These discoveries include Our nextgeneration CTLA AGEN our differentiated CD agonist AGEN our firstinclass Tregdepleting  BISPECIFIC ANTIBODY  AGEN and of course GS a bifunctional molecule now exclusively licensed to Gilead and being developed by them In summary this year we generated
Cloudera Inc	4/2019	combined company road map which we rolled out in March of this year During this period of uncertainty we saw increased competition from the  PUBLIC CLOUD  vendors Second the announcement in March of Cloudera Data Platform our new hybrid and multicloud offering created significant excitement within our customer base CDP
NVIDIA Corp	7/2015	lot of very exciting development And were working with a lot of them because we have a platform that was really designed to fuse  COMPUTER VISION  cameras from all around the car As well as radars and LIDARS and sonars and be able to do path planning and all of
Proto Labs Inc	1/2015	orders in addition we added capacity to our manufacturing facility in europe in we completed our first acquisition purchasing fineline an  ADDITIVE MANUFACTURING  or  3D PRINTING  company based in raleigh north carolina the acquisition was completed last april and is highly complementary to proto labs roughly of our customers use
Collectar Biosciences Inc	10/2017	collaboration with Acunova Therapeutics each provide these types of strategic benefits Avicenna provides us with the unique opportunity to collaborate with experts in the antibody  DRUG CONJUGATE  or ADC field Not only does this provide the opportunity to work with a very promising small molecule payload but it also allows
L-3 Communications Holdings Inc	10/2002	metal detectors where they always make you take your shoes off This is a passive scanner as I told some of you It uses  MILLIMETER WAVE  It is nonintrusive and causes no harm or disease It will guarantee you won't have a weapon on you of any kind or be

Oasis Petroleum Inc	1/2011	tell you is that the build in the backlog is really a function of the weather that we experienced and it is always difficult  FRACKING  wells in the winter but this year was particularly brutal So I think the build in the backlog was largely around the weather And then
InvenSense Inc	7/2016	as they strive to enable improved locationbased services and mapping user experience A significant opportunity for increasing our mobile content is UltraPrint our ultrasonic  FINGERPRINT SENSOR  I am very pleased to report that we are on track with the development of this gamechanging technology and have successfully passed several technology
Tesla Inc	4/2011	with our store opening in Santana Row in San Jose in April The goal here is really to engage and inform potential customers about  ELECTRIC VEHICLES  in general and the advantages of Tesla in particular and really to try to catch people before they have actually made a buying decision
SunPower Corp	10/2006	then be able to participate in the global electricity market which is measured in the form of trillion We have direct control over the  SOLAR CELL  and  SOLAR PANEL  portions of the value chain the technology core of the value chain that represents to of total installed costs In these
Vocus Inc	1/2011	content distribution along with our expansion into  SOCIAL MEDIA  Vocus is uniquely positioned to help organizations of all sizes reach and influence buyers across  SOCIAL NETWORKS  online and through the media While PR will remain a core element of the Vocus product suite we believe there is a new and
Donnelley Financial Solutions Inc	4/2018	speed and improve both the quality and consistency of business results for our clients In capital markets through the introduction of  MACHINE LEARNING  and  ARTIFICIAL INTELLIGENCE  we will improve the efficiency of XBRL tagging and align with the efforts at the SEC to move from documents to data This investment
Millennial Media LLC	4/2013	how We recently released our new Software Development Kit or SDK which is designed to enhance monetization of apps across  SMARTPHONES TABLETS  and other  CONNECTED DEVICES  SDK enhances our video advertising and rich media capabilities while adding new functionality like interactive voice ads and integration with iOSs Passbook for coupon

---

Notes: This tables presents 15 earning calls excerpts (in Column 3) with 25 words before and after technology keyword mention, with the firm (in Column 1), and the date of the earnings call (in Column 2).

Appendix Table 7 – Job postings human audit results

Panel A: Audit Results			
Audit	Use	Produce	Total
Describes company	6%	10%	16%
Describes Task	46%	34%	80%
Neither	NA	NA	4%
Panel B: Audit Results after clipping top 50 and bottom 50 words			
Audit	Use	Produce	Total
Describes company	2%	2%	4%
Describes Task	55%	36%	91%
Neither	NA	NA	5%
Panel C: Examples Excerpts			
Describes company - Produce	“[Company’s] systems offer a unique combination of technology linking <b>RFID tags</b> and sensors with displays which permit users to track locate and observe movement of equipment and people in real time currently locates millions of patient’s staff visitors and assets in healthcare facilities all over the world.”		
Describes task - Use	“passion for learning about new technology including low power RF technologies voice command systems motion control and <b>capacitive touch</b> ability to learn other non-electrical related topics mechanical and design considerations”		
Neither	“our super cool office space which doesn’t feel like an office is designed with our employees in mind techy surroundings a great outdoor space with Wi-Fi hookups for your laptop plus Bluetooth capabilities for <b>music streaming</b> we enjoy cultivating a supportive and all around positive culture that keeps our employees happy this will be a place you will want to come to everyday”		

Notes: This table presents the results from a human audit of Burning Glass technology job postings. As a part of the human audit, we classify each of randomly sampled 1,000 job postings into two types of categories 1) whether the technology keyword describes the company in the job posting or the task content of the job posting, 2) whether the job describes the use or the production of the technology. See the main text for details. In Panel A, we perform the audit on the original text of job postings for a smaller set of 100 postings. In Panel B, we clip the text of job postings by 50 words at the top and bottom, resample 1,000 postings, and then repeat the audit.

Appendix Table 8 – Job postings for technical and non-technical bigrams

Statistic	Technical bigrams (Supervised)	Non-technical bigrams (top 221)	Technical bigrams (Unsupervised)	Non-technical bigrams* (top 305)	Non-technical bigrams (ext)* (top 4000)
# bigrams	221	221	305	305	4000
Avg. postings/bigram	59,013	142	49,677	157	474
Bigrams w/ more than 100 postings	88.3%	10.0%	92.4%	9.2%	8.1%

Notes: The table presents summary statistics (the number of bigrams, the average job postings per bigram, and bigrams appearing in more than 100 job postings) for our list of 221 supervised bigrams for 29 technologies (in Column 2), the top most frequent 221 non-technical bigrams in earnings calls (in Column 3), our 305 unsupervised technical bigrams (in Column 4), the top 305 most frequent non-technical bigrams in earnings calls (in Column 5), and the top 4000 non-technical bigrams from earnings calls (in Column 6). Technical bigrams are as described in Section 2; we get to the list by intersecting bigrams in patents with bigrams in earnings calls. Non-technical bigrams are ones in earnings calls but not in patents. For both sets of bigrams, we restrict the sample to bigrams for which the share increases in earnings calls (2002-2019). \*Through the aforementioned process, we obtained many more non-technical (104,627) bigrams than supervised bigrams (221) and technical bigrams (305). We restrict the sample to the top (by frequency in earnings calls) 221 (in Column 3), 305 non-technical bigrams (in Column 5) and 4,000 non-technical bigrams (in Column 6).

Appendix Table 9 – Top technical and non-technical bigrams

Top technical bigrams			Top non-technical bigrams		
Bigram	# earnings calls	# job postings	Bigram	# earnings calls	# job postings
mobile devices	6597	1078049	bofa merrill	34490	221
machine learning	2860	525286	stifel nicolaus	28877	256
cloud computing	2781	485333	division associate	12472	4237
cloud services	2450	380980	keefe bruyette	11682	16
quality metrics	2029	196497	bruyette woods	11498	14

Notes: The table presents the top five most frequent technical and non-technical bigrams in earnings calls. Technical bigrams are as described in Section 2; we get the list by intersecting bigrams in patents with bigrams in earnings calls. Non-technical bigrams are ones in earnings calls but not in patents. For both sets of bigrams, we restrict to the sample to bigrams for which the share increases in earnings calls (2002-2019).

Appendix Table 10 – Top occupations by technology exposure (Virtual Reality)

Occupations	Normalized Share	Total Job Postings
Computer Hardware Engineers	36.56	8,486
Fine Artists, Including Painters, Sculptors, and Illustrators	30.77	6,147
Computer and Information Research Scientists	25.03	23,241
Multimedia Artists and Animators	23.96	6,461
Art Directors	23.69	7,641
Computer Science Teachers, Postsecondary	16.66	3,626
Communications Teachers, Postsecondary	16.32	1,994
Aerospace Engineering and Operations Technicians	15.9	633
Sound Engineering Technicians	14.9	2,804
Social Science Research Assistants	13.92	5,605
Biomedical Engineers	12.86	1,848
Aircraft Mechanics and Service Technicians	10.75	11,360
Producers and Directors	9.89	14,838
Models	9.83	1,733
Commercial and Industrial Designers	8.3	18,572
Psychology Teachers, Postsecondary	7.96	3,430
Interior Designers	7.77	9,449
Health Specialties Teachers, Postsecondary	7.63	7,508
Natural Sciences Managers	7.53	31,612
Art, Drama, and Music Teachers, Postsecondary	6.8	2,415

Notes: This table lists the top occupations (in Column 1), their normalized share of technology postings (in Column 2), and the total job postings for the occupation (in Column 3) associated with “Virtual Reality.” The normalized share and total job postings in the table are calculated after excluding the years before the year of emergence of the technology.

Appendix Table 11 – Top occupations by technology exposure (all technologies)

Technology	Top Exposed Occupations (by normalized share)
3D printing	Materials Engineers(44.52);Materials Scientists(37.40);
Autonomous cars	Computer Hardware Engineers(39.15);Computer and Information Research Scientists(23.25);
Bispecific monoclonal antibody	Biological Technicians(51.82);Biological Scientists, All Other(46.54);
Cloud computing	Sales Engineers( 9.98);Computer Network Architects( 8.37);
Computer vision	Computer and Information Research Scientists(53.19);Computer Hardware Engineers(46.67);
Drug conjugates	Chemical Technicians(53.19);Biological Scientists, All Other(53.19);
Electronic gaming	Fine Artists, Including Painters, Sculptors, and Illustrators(46.59);Gaming Service Workers, All Other(39.11);
Extremely high frequency	Electronics Engineers, Except Computer(53.19);Computer Hardware Engineers(47.63);
Fingerprint sensor	Computer Hardware Engineers(41.31);Human Resources Assistants, Except Payroll and Timekeeping(36.20);
Fracking	Petroleum Engineers(53.19);Geoscientists, Except Hydrologists and Geographers(39.40);
GPS	Surveyors(53.19);Surveying and Mapping Technicians(53.19);
Hybrid vehicle/Electric car	Power Plant Operators(26.60);Solar Photovoltaic Installers(26.05);
Lane departure warning	Mechanical Engineers(22.66);Engineers, All Other(14.74);
Lithium battery	Materials Scientists(43.52);Materials Engineers(40.26);
Machine learning/AI	Computer and Information Research Scientists(53.19);Computer Science Teachers, Postsecondary(14.12);
Mobile payment	Food Scientists and Technologists(18.59);Marketing Managers(10.57);
OLED display	Materials Scientists(21.77);Computer Hardware Engineers(20.76);
Online streaming	Audio and Video Equipment Technicians(43.63);Film and Video Editors(38.82);
RFID	Locksmiths and Safe Repairers(20.75);Electronics Engineers, Except Computer(17.69);
Search engine	Writers and Authors(17.93);Advertising and Promotions Managers(15.31);
Smart devices	Electronic Equipment Installers and Repairers, Motor Vehicles(13.08);Automotive Glass Installers and Repairers(10.04);
Social networking	Reporters and Correspondents(20.26);Public Relations Specialists(16.62);
Software defined radio	Electronics Engineers, Except Computer(52.28);Computer Hardware Engineers(49.58);
Solar power	Solar Photovoltaic Installers(53.19);Wind Turbine Service Technicians(24.56);
Stent graft	Sales Representatives (Technical and Scientific Products) (38.01);Cardiovascular Technologists (34.56);
Touch screen	Audio and Video Equipment Technicians(28.19);Multimedia Artists and Animators(14.45);
Virtual reality	Computer Hardware Engineers(36.56);Fine Artists, Including Painters, Sculptors, and Illustrators(30.77);
Wi-Fi	Electronic Home Entertainment Equipment Installers and Repairers(40.58);Electronics Engineers, Except Computer(23.64);
Wireless charging	Computer Hardware Engineers(51.07);Electrical Engineers(42.92);

Notes: This table lists the top exposed occupations (in Column 2) for each of our 29 technologies (in Column 1), and the normalized share of postings exposed to the technology (in parentheses alongside each occupation).



Appendix Table 12 – Top regions by technology exposure (all technologies)

Technology	Top Exposed CBSAs (by normalized share)
3D printing	Los Alamos, NM (10.63);Corvallis, OR ( 9.87);
Autonomous cars	San Jose-Sunnyvale-Santa Clara, CA (13.66);Detroit-Warren-Dearborn, MI (13.47);
Bispecific monoclonal antibody	Worcester, MA-CT ( 8.80);San Francisco-Oakland-Hayward, CA ( 8.26);
Cloud computing	San Jose-Sunnyvale-Santa Clara, CA ( 4.95);San Francisco-Oakland-Hayward, CA ( 3.28);
Computer vision	San Jose-Sunnyvale-Santa Clara, CA (10.27);Trenton, NJ ( 6.70);
Drug conjugates	Seattle-Tacoma-Bellevue, WA ( 6.90);San Francisco-Oakland-Hayward, CA ( 6.58);
Electronic gaming	Seattle-Tacoma-Bellevue, WA ( 6.73);Reno, NV ( 4.84);
Extremely high frequency	Atlantic City-Hammonton, NJ (13.08);Manchester-Nashua, NH (11.23);
Fingerprint sensor	San Jose-Sunnyvale-Santa Clara, CA (11.97);Buffalo-Cheektowaga-Niagara Falls, NY ( 6.87);
Fracking	Williston, ND (17.33);Midland, TX (17.33);
GPS	Butte-Silver Bow, MT ( 9.37);Warner Robins, GA ( 8.91);
Hybrid vehicle/Electric car	Milwaukee-Waukesha-West Allis, WI ( 8.85);Detroit-Warren-Dearborn, MI ( 7.90);
Lane departure warning	Detroit-Warren-Dearborn, MI (11.39);Ann Arbor, MI ( 6.38);
Lithium battery	Midland, MI (16.64);Joplin, MO (15.81);
Machine learning/AI	San Jose-Sunnyvale-Santa Clara, CA ( 7.47);Los Alamos, NM ( 7.44);
Mobile payment	San Jose-Sunnyvale-Santa Clara, CA ( 5.09);Newton, IA ( 4.41);
OLED display	San Jose-Sunnyvale-Santa Clara, CA (13.63);Corning, NY ( 6.50);
Online streaming	Athens, TX ( 3.81);Green Bay, WI ( 3.04);
RFID	Fond du Lac, WI ( 7.60);Bellefontaine, OH ( 3.90);
Search engine	Orangeburg, SC ( 3.55);Oxford, NC ( 2.94);
Smart devices	San Jose-Sunnyvale-Santa Clara, CA ( 3.84);Scottsbluff, NE ( 3.06);
Social networking	Ellensburg, WA ( 3.41);Truckee-Grass Valley, CA ( 3.00);
Software defined radio	Cedar Rapids, IA (16.05);Rochester, NY (12.22);
Solar power	San Luis Obispo-Paso Robles-Arroyo Grande, CA ( 6.82);Toledo, OH ( 6.78);
Stent graft	Santa Rosa, CA (18.17);Tampa-St. Petersburg-Clearwater, FL (12.50);
Touch screen	Edwards, CO ( 9.09);Breckenridge, CO ( 8.26);
Virtual reality	Marshalltown, IA (11.96);Hinesville, GA ( 8.33);
Wi-Fi	Berlin, NH-VT ( 5.33);Athens, TX ( 4.36);
Wireless charging	San Jose-Sunnyvale-Santa Clara, CA (13.41);Youngstown-Warren-Boardman, OH-PA (10.09);

Notes: This table lists the top exposed Core-Based Statistical Areas (in Column 2) for each of our 29 technologies (in Column 1), and the normalized share of postings exposed to the technology (in parentheses alongside each CBSA).

Appendix Table 13 – Top industries by technology exposure (all technologies)

Technology	Top Exposed Industries (by normalized share)
3D printing	Metalworking Machinery Manufacturing(17.44);Specialized Design Services(16.55);
Autonomous cars	Household Appliance Manufacturing(52.51);Motor Vehicle Parts Manufacturing(47.50);
Bispecific monoclonal antibody	Pharmaceutical and Medicine Manufacturing(53.77);Scientific Research and Development Services(20.14);
Cloud computing	Computer and Peripheral Equipment Manufacturing(14.76);Software Publishers(12.86);
Computer vision	Space Research and Technology (29.37);Semiconductor and Other Electronic Component Manufacturing(25.07);
Drug conjugates	Pharmaceutical and Medicine Manufacturing(48.77);Scientific Research and Development Services(29.37);
Electronic gaming	Motion Picture and Video Industries(20.48);Electronic Shopping and Mail-Order Houses (19.47);
Extremely high frequency	Satellite Telecommunications (41.25); Navigational, Measuring, Electromedical .. Manufacturing(40.43);
Fingerprint sensor	Semiconductor and Other Electronic Component Manufacturing(35.78);Junior Colleges(17.08);
Fracking	Oil and Gas Extraction(60.02);Support Activities for Mining(60.02);
GPS	Grocery and Related Product Merchant Wholesalers (46.28);Support Activities for Forestry(39.57);
Hybrid vehicle/Electric car	Motor Vehicle Manufacturing(45.86);Electric Power Generation, Transmission and Distribution (43.00);
Lane departure warning	Motor Vehicle Parts Manufacturing(44.42);Household Appliance Manufacturing(24.01);
Lithium battery	Other Electrical Equipment and Component Manufacturing(54.57);Plastics Product Manufacturing(31.80);
Machine learning/AI	Electronic Shopping and Mail-Order Houses (19.03);Other Information Services(13.30);
Mobile payment	Activities Related to Credit Intermediation (43.72);Electronic Shopping and Mail-Order Houses (23.04);
OLED display	Audio and Video Equipment Manufacturing(37.30);Semiconductor and Other Electronic Component Manufacturing(19.30);
Online streaming	Motion Picture and Video Industries(19.69);Radio and Television Broadcasting(19.57);
RFID	Industrial Machinery Manufacturing(42.11);Converted Paper Product Manufacturing(21.87);
Search engine	Other Information Services(53.81);Newspaper, Periodical, Book, and Directory Publishers( 8.26);
Smart devices	Electronics and Appliance Stores (19.77);Audio and Video Equipment Manufacturing(13.41);
Social networking	Motion Picture and Video Industries(10.91);Radio and Television Broadcasting( 9.93);
Software defined radio	Communications Equipment Manufacturing(40.49);Aerospace Product and Parts Manufacturing(39.00);
Solar power	Electric Power Generation, Transmission and Distribution (55.57);Hardware, and Plumbing and Heating Equipment Wholesalers (54.57);
Stent graft	Navigational, Measuring, Electromedical ..Instruments Manufacturing (52.42); Resin, .. and Artificial Fibers and Filaments Manfcn(16.37);
Touch screen	Printing and Related Support Activities (14.51); Resin, Synthetic Rubber, and Artificial .. Fibers and Filaments Manufacturing(13.82);
Virtual reality	Other Information Services(25.88);Audio and Video Equipment Manufacturing(22.77);
Wi-Fi	Wireless Telecommunications Carriers(21.63);Cable and Other Subscription Programming(18.68);
Wireless charging	Semiconductor and Other Electronic Component Manufacturing(59.42);Motor Vehicle Parts Manufacturing(20.20);

Notes: This table lists the top exposed industries (in Column 2) for each of our 29 technologies (in Column 1), and the normalized share of postings exposed to the technology (in parentheses alongside each industry).

Appendix Table 14 – Year of emergence by technology

Technology	Year of Emergence	
	EC (baseline)	Patents
3D printing	2011	2013
Autonomous cars	2014	2012
Bispecific antibody	2012	1999
Cloud computing	2008	2011
Computer vision	2008	2006
Drug conjugates	2002*	2002
Electronic gaming	2002*	1995
Millimeter wave	2014	2012
Fingerprint sensor	2011	2005
Fracking	2007	2005
GPS	2002*	1999
Hybrid vehicle/Electric car	2007	2006
Lane departure warning	2002*	2004
Lithium battery	2002*	1994
Machine learning/AI	2015	2005
Mobile payment	2007	2007
OLED display	2002*	2005
Online streaming	2002*	1997
RFID	2002*	2004
Search engine	2002*	1997
Smart devices	2005	2010
Social networking	2006	2009
Software defined radio	2002*	2005
Solar power	2002*	1975
Stent graft	2003	1995
Touch screen	2002*	2010
Virtual reality	2013	2012
Wi-Fi	2002*	2007
Wireless charging	2012	2012

Notes: The table provides our set of technologies (in Column 1), their year of emergence based on our earnings call data (Column 2), and their alternative emergence year based on patent data (Column 3). The year of emergence in Column 2 is calculated as the first year that the share of firms mentioning the technology in earnings calls reaches 10% of its maximum between 2002 and 2019. Years of emergence marked with \* are technologies for which share of firms mentioning them in 2002 is already more than 10% of the maximum share of firms mentioning them over the sample period 2002-19. For these technologies, we impute the emergence year to be 2002. In Column 3, the year of emergence is the year in which the share of U.S. patents applied for in a technology reaches 50% of its maximum value between 1976 and 2015.

Appendix Table 15 - Summary statistics

	Mean	SD	p25	p50	p75	N
	(1)	(2)	(3)	(4)	(5)	(6)
Panel A: Location						
Normalized Share	0.533	3.801	0	0	0.301	266467
University assets per capita	5670.065	12314.42	411.181	2234.962	5455.663	917
University enrollment per capita	0.141	0.178	0.024	0.1	0.171	917
Share of college educated	0.198	0.065	0.148	0.182	0.237	917
Share of post-graduate	0.068	0.028	0.049	0.06	0.081	917
Coefficient of Variation	3.716	2.471	1.729	3.033	5.181	249
Panel B: Industry						
Normalized Share	1.4	18.267	0	0	0.253	88490
Coefficient of Variation	4.885	2.315	3.433	4.429	5.867	249
Panel C: Occupation						
Normalized Share	0.704	3.851	0	0	0.075	238777
Share College Educated	54.774	13.125	47.462	56.785	63.325	249
Share Post Graduate	19.098	7.054	14.825	18.18	22.309	249
Wage	63044.66	11143.04	57776.79	64014.05	71611.97	249
Years of Schooling	14.993	0.769	14.579	15.064	15.403	249
Coefficient of Variation	6.654	3.32	3.914	5.853	8.909	249
Panel D: Firms						
Normalized Share	0.591	14.258	0	0	0	38990627
Coefficient of Variation	34.515	25.759	16.342	29.253	44.908	249

Notes: This tables shows summary statistics for variables used in the analyses of the paper. Summary statistics (Columns 2-6) are shown for the pooled sample of technologies including and after the year of emergence. In our sample, location is one of 917 Core-Based Statistical Areas (CBSA), industry is one of 311 4-digit North American Industry Classification System (NAICS) codes, occupation is one of 836 six-digit Standard Occupation Classification (SOC) codes, and a firm is one of 329,158 unique firms in BG. The normalized share of technology jobs in all panels is calculated as  $Normalized\ share_{i,\tau,t} = \frac{share\ jobs\ exposed_{i,\tau,t}}{share\ jobs\ exposed_{\tau,t}}$ , where  $i$  is a location, industry, occupation, or firm (cell). The coefficient of variation in all panels is calculated over normalized share of technology job postings over cells for technology x year observations. To calculate the coefficient of variation across firms in BG, we keep the sample across years comparable by keeping only 10,231 firms which post at least one job postings in each year in the BG sample. Location variables (in Panel A, Rows 2-5) are reported in the table for the cross-section of 917 CBSAs and calculated as following: university assets per capita is calculated as the total assets reported by universities in a CBSA in the Higher Education Research and Development (HERD) survey and normalized by the population of the CBSA; and enrollment per capita is calculated as the total enrollment reported by universities in a CBSA in the HERD survey and normalized by the population of the CBSA. Skill level variables (in Panel C, Rows 2-5) are calculated using  $Skill_t^{\tau} = \frac{\sum_o N_{o,t}^{\tau} \chi_{o;2015}}{\sum_o N_{o,t}^{\tau}}$ , where  $\chi_{o;2015}$  is the skill measure of interest (eg: share of college-educated people) in an occupation (o) in the 2015 American Community Survey and  $N_{o,t}^{\tau}$  is the number of technology job postings in technology  $\tau$ , occupation o, and time t.

Appendix Table 16 – Top Pioneer locations by technology

Technology	Top CBSA Pioneer	State
3D printing	Boston-Cambridge-Newton	MA-NH
Autonomous cars	San Jose-Sunnyvale-Santa Clara	CA
Bispecific antibody	San Francisco-Oakland-Hayward	CA
Cloud computing	San Jose-Sunnyvale-Santa Clara	CA
Computer vision	San Jose-Sunnyvale-Santa Clara	CA
Drug conjugates	Boston-Cambridge-Newton	MA-NH
Electronic gaming	San Jose-Sunnyvale-Santa Clara	CA
Millimeter wave	New York-Newark-Jersey City	NY-NJ-PA
Fingerprint sensor	San Jose-Sunnyvale-Santa Clara	CA
Fracking	Houston-The Woodlands-Sugar Land	TX
GPS	San Jose-Sunnyvale-Santa Clara	CA
Hybrid vehicle/Electric car	Detroit-Warren-Dearborn	MI
Lane departure warning	Grand Rapids-Wyoming	MI
Lithium battery	Los Angeles-Long Beach-Anaheim	CA
Machine learning/AI	San Jose-Sunnyvale-Santa Clara	CA
Mobile payment	San Francisco-Oakland-Hayward	CA
OLED display	Trenton	NJ
Online streaming	San Jose-Sunnyvale-Santa Clara	CA
RFID	Grand Rapids-Wyoming	MI
Search engine	San Jose-Sunnyvale-Santa Clara	CA
Smart devices	San Jose-Sunnyvale-Santa Clara	CA
Social networking	San Jose-Sunnyvale-Santa Clara	CA
Software defined radio	Boulder	CO
Solar power	San Jose-Sunnyvale-Santa Clara	CA
Stent graft	San Francisco-Oakland-Hayward	CA
Touch screen	San Jose-Sunnyvale-Santa Clara	CA
Virtual reality	San Jose-Sunnyvale-Santa Clara	CA
Wi-Fi	New York-Newark-Jersey City	NY-NJ-PA
Wireless charging	Boston-Cambridge-Newton	MA-NH

Notes: This table shows the top Pioneer location (in Column 2) for each of our 29 technologies (Column 1), and its state(s) (Column 3). We define Pioneer locations as those which collectively account for 50% of patent grants associated with a given technology and applied for within ten years before its emergence. The top Pioneer location listed in the table is the one with most patents.

Appendix Table 17 – Region broadening by skill

	Coefficient of Variation across Locations		
	(1)	(2)	(3)
	Low Skill	Medium Skill	High Skill
<i>Years since emergence</i> $_{\tau,t}$	-0.154*** (0.049)	-0.169*** (0.048)	-0.097*** (0.033)
R2	0.841	0.851	0.916
N	231	231	231

Notes: This table reports the results from regressions of coefficient of variation during lifecycle of a technology on year since emergence of the technology, separately for low skill occupations (Column 1), medium skill occupations (Column 2), and high skill occupations (Column 3). To calculate the coefficient of variation by skill, we aggregate the job postings data over occupation, CBSA, and year, and then separately for high-skill occupations (with the share of college educated people > 60%), medium-skill occupations (with the share of college educated people > 30% and <60%), and low-skill occupations (with the share of college educated people < 30%). The coefficient of variation is calculated over *Normalized share* $_{cbsa,\tau,t,skill}$  across CBSAs by skill group, technology, and time. All specifications control for technology and year fixed effects. Standard errors are clustered by technology.

Appendix Table 18 – Pioneer occupations and industries by technology

Technology	Top Pioneer Occupation (share of patents)	Top Pioneer Industry (share of patents)
3D printing	Mechanical Engineers (0.140)	Computer and Peripheral Equipment Manufacturing (0.419)
Autonomous cars	Computer Occupations All Other (0.186)	Motor Vehicle Manufacturing (0.370)
Bispecific monoclonal antibody	Operations Research Analysts (0.375)	Pharmaceutical and Medicine Manufacturing (0.946)
Cloud computing	Software Developers Applications (0.228)	Software Publishers (0.300)
Computer vision	Software Developers Applications (0.295)	Semiconductor and Other Electronic Component Manufacturing (0.174)
Drug conjugates	Natural Sciences Managers (0.135)	Pharmaceutical and Medicine Manufacturing (0.910)
Electronic gaming	Software Developers Applications (0.202)	Software Publishers (0.202)
Millimeter wave	Electronics Engineers Except Computer (0.169)	Semiconductor and Other Electronic Component Manufacturing (0.371)
Fingerprint sensor	Software Developers Applications (0.203)	Semiconductor and Other Electronic Component Manufacturing (0.215)
Fracking	Geoscientists Except Hydrologists and Geographers (0.286)	Oil and Gas Extraction (0.881)
GPS	Computer Occupations All Other (0.173)	Communications Equipment Manufacturing (0.187)
Hybrid vehicle/Electric car	Mechanical Engineers (0.151)	Motor Vehicle Manufacturing (0.681)
Lane departure warning	Mechanical Engineers (0.500)	Motor Vehicle Manufacturing (0.393)
Lithium battery	Electrical Engineers (0.188)	Commercial and Service Industry Machinery Manufacturing (0.115)
Machine learning/AI	Software Developers Applications (0.251)	Other Information Services (0.225)
Mobile payment	Marketing Managers (0.154)	Semiconductor and Other Electronic Component Manufacturing (0.227)
OLED display	Engineers All Other (0.400)	Commercial and Service Industry Machinery Manufacturing (0.320)
Online streaming	Sales Representatives (0.095)	Semiconductor and Other Electronic Component Manufacturing (0.188)
RFID tags	Architectural and Engineering Managers (0.098)	Computer and Peripheral Equipment Manufacturing (0.191)
Search engine	Marketing Managers (0.124)	Other Information Services (0.264)
Smart devices	Software Developers Applications (0.229)	Software Publishers (0.243)
Social networking	Marketing Managers (0.128)	Other Information Services (0.299)
Software defined radio	Software Developers Applications (0.489)	Communications Equipment Manufacturing (0.353)
Solar power	Mechanical Engineers (0.099)	Semiconductor and Other Electronic Component Manufacturing (0.243)
Stent graft	Physicians and Surgeons All Other (0.375)	Medical Equipment and Supplies Manufacturing (0.628)
Touch screen	Sales Representatives Wholesale (0.134)	Commercial and Service Industry Machinery Manufacturing (0.211)
Virtual reality	Software Developers Applications (0.198)	Semiconductor and Other Electronic Component Manufacturing (0.214)
Wi-Fi	Retail Salespersons (0.255)	Communications Equipment Manufacturing (0.314)
Wireless charging	Computer Occupations All Other (0.222)	Semiconductor and Other Electronic Component Manufacturing (0.412)

Notes: The table shows the top Pioneer occupation (Column 2) and top Pioneer industry (Column 3) for each of our 29 technologies (in Column 1). A Pioneer is defined as the set of occupations and industries that account for more than 50% of patents associated with the technology during the ten years before year of emergence of the technology. The top Pioneer is the one with most patents during the same period. See Section 4 of the main text for details.

Appendix Table 19 – Robustness: Emergence year, job postings sample, and std. errors

Panel A: Patent Emergence Year				
	(1) Region Broadening	(2) Pioneer Persistence	(3) Skill Broadening	(4) Region Broadening by Skill
<i>Years since emergence<sub>τ,t</sub></i> ( <i>patents</i> )	-0.070*** (0.020)		-0.727*** (0.226)	-0.121*** (0.040)
<i>Pioneer</i>		1.369*** (0.410)		
<i>Pioneer</i> * <i>yrs since emg<sub>τ,t</sub>(patents)</i>		-0.033** (0.014)		
<i>Years since emergence<sub>τ,t</sub></i> ( <i>patents</i> ) * <i>1{skill = low}</i>				-0.046* (0.026)
R2	0.893	0.077	0.880	0.750
N	255	275,751	255	510
Panel B: Without 2007				
	(1)	(2)	(3)	(4)
<i>Years since emergence<sub>τ,t</sub></i>	-0.100*** (0.036)		-0.877*** (0.272)	-0.089* (0.046)
<i>Pioneer</i>		2.475*** (0.643)		
<i>Pioneer</i> * <i>yrs since emg<sub>τ,t</sub></i>		-0.157*** (0.048)		
<i>Years since emg<sub>τ,t</sub></i> * <i>1{skill = low}</i>				-0.184*** (0.044)
R2	0.891	0.079	0.880	0.780
N	236	248,873	236	504
Panel C: Robust Standard Errors				
	(1)	(2)	(3)	(4)
<i>Years since emergence<sub>τ,t</sub></i>	-0.092*** (0.023)		-0.919*** (0.224)	-0.110*** (0.027)
<i>Pioneer</i>		2.313*** (0.202)		
<i>Pioneer</i> * <i>yrs since emg<sub>τ,t</sub></i>		-0.146*** (0.016)		
<i>Years since emg<sub>τ,t</sub></i> * <i>1{skill = low}</i>				-0.195*** (0.023)
R2	0.888	0.075	0.873	0.772
N	249	266,467	249	538

Notes: This table reports results of robustness checks for our primary results. The four columns replicate our baseline specifications from Table 4, Column 1; Table 5, Column 2; Table 6 Column 1; and Table 7, Column 3; respectively. In Panel A, we calculate the year of emergence as the year in which the share of US patents for a technology reaches 50% of their maximum value between 1976 and 2015; in Panel B, we exclude the year 2007; in Panel C, we use robust standard errors instead of the clustered ones in the baseline specification. In all panels, Column 1, 3, and 4 control for technology and year fixed effects, while column 2 controls for CBSA, technology and year fixed effects. Standard errors are clustered in Panels A and B; they are robust in Panel C.



Appendix Table 20 - Robustness: Skill broadening with sample reweighted to U.S. employment

	(1) Share of college educated * 100	(2) Share of post graduate * 100	(3) Avg. wage	(4) Avg. schooling
<i>Years since emergence<sub>t,t</sub></i>	-0.593** (0.272)	-0.180 (0.118)	-627.958** (241.790)	-0.035** (0.015)
R2	0.902	0.915	0.907	0.905
N	249	249	249	249

Notes: This table explores the robustness of our skill-broadening result. We regress the approximate skill composition of technology job postings,  $Skill_t^T$ , on the years since the inception of the technology. In this case, the technology postings in an occupation are reweighted to match the composition of hiring in the U.S. economy for each two-digit (SOC) occupation as described in Appendix 3.3. Otherwise all specifications are identical to those in Table 6.

Appendix Table 21 – Robustness: Alternative cutoffs defining disruptive technologies

	(1)	(2)	(3)	(4)	(5)
	$\frac{EC_{2002}}{EC_{max}} \leq 0.01$	$\frac{EC_{2002}}{EC_{max}} \leq 0.05$	$\frac{EC_{2002}}{EC_{max}} \leq 0.1$	$\frac{EC_{2002}}{EC_{max}} \leq 0.15$	$\frac{EC_{2002}}{EC_{max}} \leq 0.2$
Panel A: Skill broadening					
	Share College Educated				
	(1)	(2)	(3)	(4)	(5)
<i>Years since emergence<sub>τ,t</sub></i>	--0.315*** (0.117)	--0.329*** (0.113)	--0.363*** (0.095)	--0.334*** (0.073)	--0.312*** (0.061)
R2	0.888	0.882	0.873	0.875	0.877
N	1,740	1,825	2,339	3,138	3,953
Panel B: Region broadening					
	Coefficient of Variation				
	(1)	(2)	(3)	(4)	(5)
<i>Years since emergence<sub>τ,t</sub></i>	--0.123*** (0.020)	--0.129*** (0.020)	--0.150*** (0.017)	--0.141*** (0.013)	--0.129*** (0.011)
	0.844	0.839	0.841	0.847	0.850
	1,740	1,825	2,339	3,138	3,953
Panel C: Pioneer location persistence					
	Normalized Share of Technology Postings				
	(1)	(2)	(3)	(4)	(5)
<i>Pioneer</i>	1.484*** (0.449)	1.436*** (0.428)	1.412*** (0.397)	1.308*** (0.370)	1.230*** (0.341)
<i>Pioneer * yrs since emg<sub>τ,t</sub></i>	--0.080** (0.033)	--0.075** (0.031)	--0.080*** (0.029)	--0.068*** (0.026)	--0.063*** (0.024)
R2	0.021	0.021	0.022	0.021	0.022
N	2,097,244	2,178,014	2,757,321	3,653,199	4,590,851
Panel D: Differential decline in concentration by skill					
	Coefficient of Variation (for low and high skill occupations)				
	(1)	(2)	(3)	(4)	(5)
<i>Years since emergence<sub>τ,t</sub></i>	--0.105*** (0.034)	--0.096*** (0.033)	--0.094*** (0.029)	--0.068*** (0.024)	--0.061*** (0.020)
<i>Years since emg<sub>τ,t</sub></i> <i>* 1{skill = low}</i>	--0.171*** (0.030)	--0.180*** (0.030)	--0.205*** (0.025)	--0.205*** (0.020)	--0.199*** (0.016)
R2	0.686	0.683	0.691	0.705	0.717
N	3,469	3,639	4,663	6,261	7,886

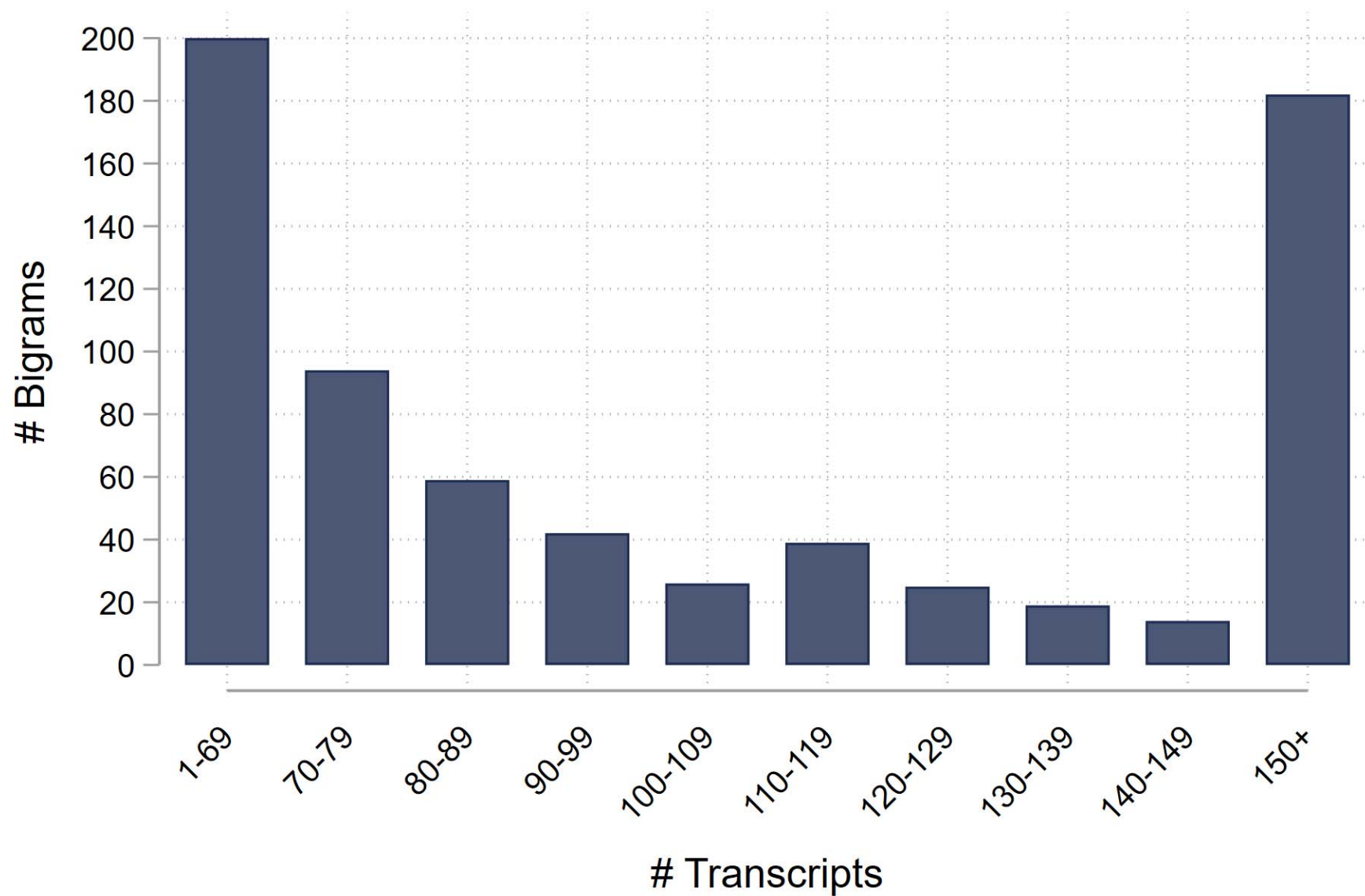
Notes: This table reports our primary results, replicated for varying disruption cut-offs as given in Appendix Table 1. In each case, we treat all (235, 243, 305, 401, and 500) bigrams as a separate technology without further human intervention, as in our unsupervised approach. In Panel A, we replicate our primary results in Table 4, Column 1; Panel B corresponds to Table 5, Column 2; Panel C corresponds to Table 6, Column 1; and Panel D corresponds to Table 7, Column 3.

Appendix Table 22 – Robustness: Alternative minimum number of mentions in earnings calls

	(1)	(2)	(3)	(4)	(5)
	# $EC \geq 80$	# $EC \geq 90$	# $EC \geq 100$	# $EC \geq 110$	# $EC \geq 120$
Panel A: Skill broadening					
	Share College Educated				
	(1)	(2)	(3)	(4)	(5)
<i>Years since emergence<sub><math>\tau,t</math></sub></i>	--0.342*** (0.075)	--0.364*** (0.087)	--0.363*** (0.095)	--0.375*** (0.103)	--0.336*** (0.107)
R2	0.884	0.874	0.873	0.867	0.866
N	3,004	2,623	2,339	2,135	1,842
Panel B: Region broadening					
	Coefficient of Variation				
	(1)	(2)	(3)	(4)	(5)
<i>Years since emergence<sub><math>\tau,t</math></sub></i>	--0.141*** (0.014)	--0.145*** (0.016)	--0.150*** (0.017)	--0.156*** (0.018)	--0.156*** (0.020)
	0.848	0.844	0.841	0.836	0.835
	3,004	2,623	2,339	2,135	1,842
Panel C: Pioneer location persistence					
	Normalized Share of Technology Postings				
	(1)	(2)	(3)	(4)	(5)
<i>Pioneer</i>	1.410*** (0.385)	1.416*** (0.405)	1.412*** (0.397)	1.393*** (0.378)	1.416*** (0.376)
<i>Pioneer * yrs since emg<sub><math>\tau,t</math></sub></i>	--0.078*** (0.027)	--0.080*** (0.029)	--0.080*** (0.029)	--0.080*** (0.027)	--0.080*** (0.026)
R2	0.022	0.022	0.022	0.023	0.023
N	3,604,929	3,109,176	2,757,321	2,502,947	2,150,159
Panel D: Differential decline in concentration by skill					
	Coefficient of variation (for low and high skill occupations)				
	(1)	(2)	(3)	(4)	(5)
<i>Years since emergence<sub><math>\tau,t</math></sub></i>	--0.105*** (0.025)	--0.100*** (0.027)	--0.094*** (0.029)	--0.096*** (0.030)	--0.100*** (0.032)
<i>Years since emg<sub><math>\tau,t</math></sub></i> * 1{skill = low}	--0.191*** (0.021)	--0.195*** (0.024)	--0.205*** (0.025)	--0.208*** (0.027)	--0.201*** (0.029)
R2	0.693	0.692	0.691	0.689	0.689
N	5,991	5,229	4,663	4,257	3,676

Notes: This table reports our primary results, replicated for varying minimum numbers of mentions in earnings calls as given in Appendix Table 2. In Panel A, we replicate our primary results in Table 4, Column 1; Panel B corresponds to Table 5, Column 2; Panel C corresponds to Table 6, Column 1; and Panel D corresponds to Table 7, Column 3.

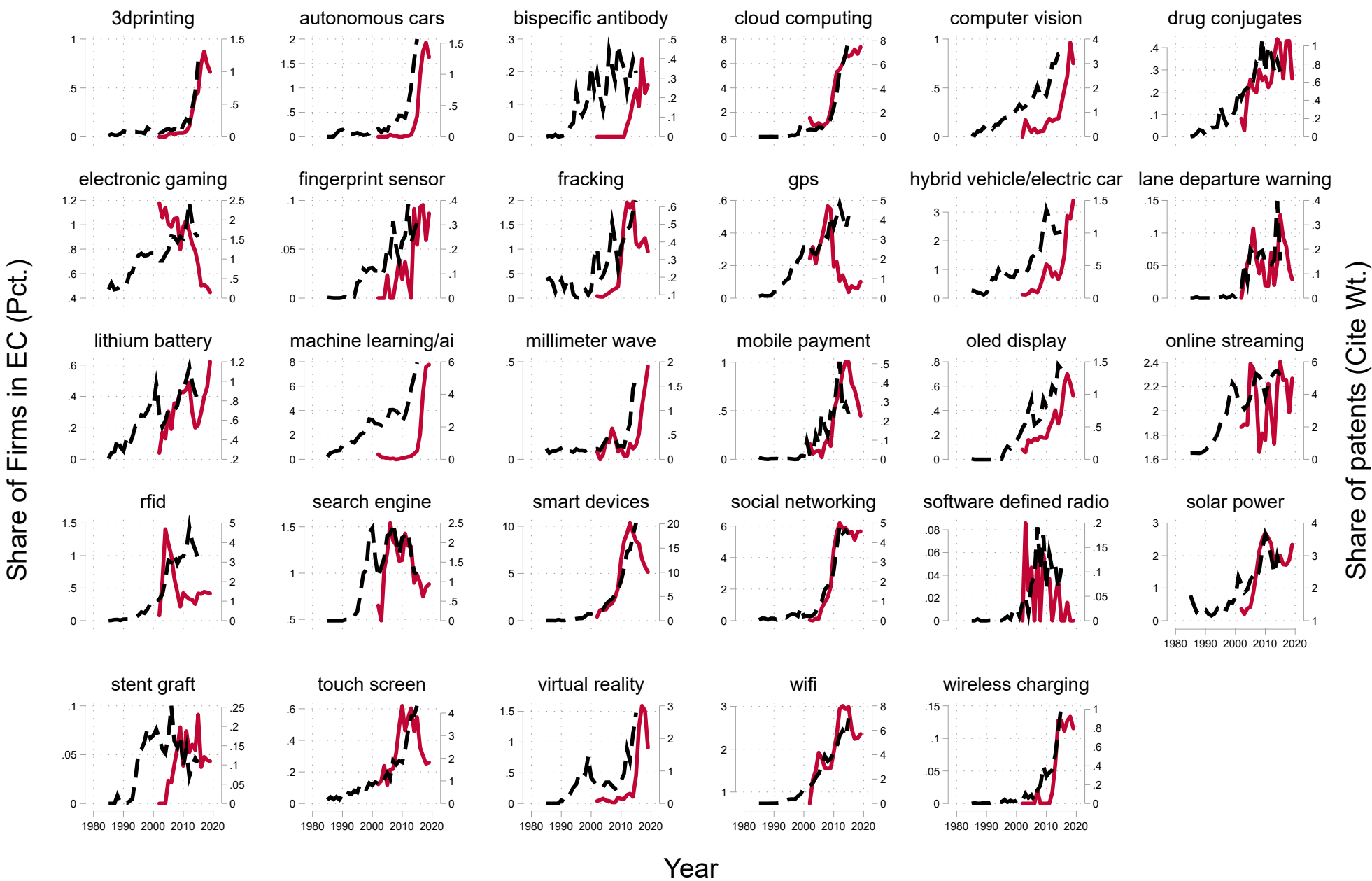
**Appendix Figure 1 – Distribution of technical bigrams**



**Notes:** This figure plots the number of technical bigrams that increase at least 10-fold in our sample of earnings calls and the number of earnings calls that they are mentioned in. We cap the first bar for the number of bigrams which appear in 1-69 earnings calls at 200 for visual clarity. A total of 13,730 bigrams appear in 1-69 earnings calls.

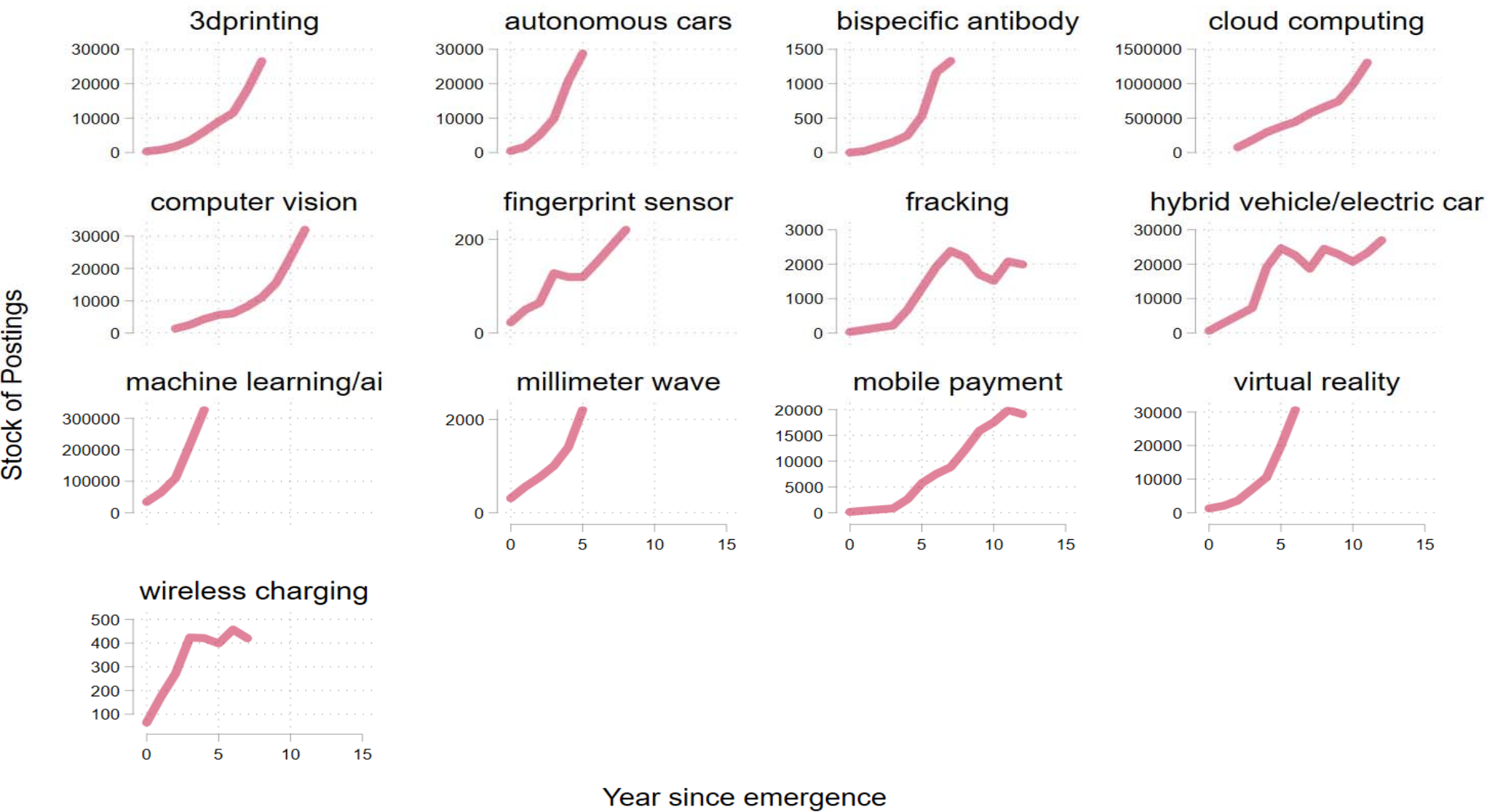


Appendix Figure 3 – Exposed patents (in black-dashed) and earnings calls (in red-solid), by technology and year



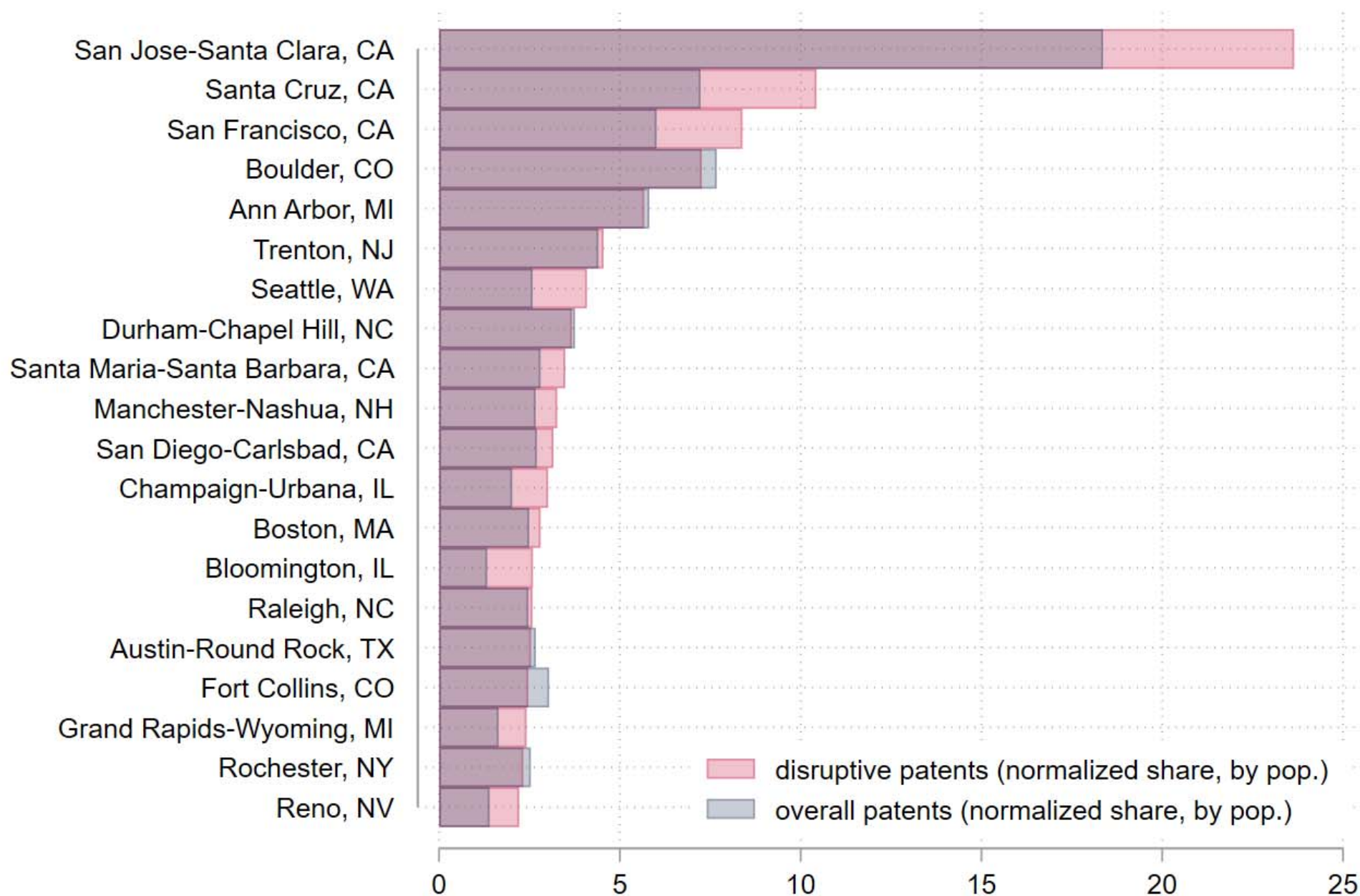
**Notes:** The figure plots the share of firms in earnings calls exposed to each technology (in red-solid), and the share of citation-weighted patents associated with each of our 29 technologies in (black-dashed). The sample includes earnings calls between 2002 and 2019 and of patents between 1985 and 2015. The overall correlation between the two time series is 80.26%.

Appendix Figure 4 – Stock of technology postings, by year since emergence



**Notes:** The figure plots stocks of technology job postings against year since emergence for 13 technologies with a year of emergence post-2007 (beginning of the job postings sample). For each technology, in any given year, the stock of technology postings is calculated by cumulatively adding technology postings and assuming a 40% separation rate. We also assume no technology postings before the year of emergence.

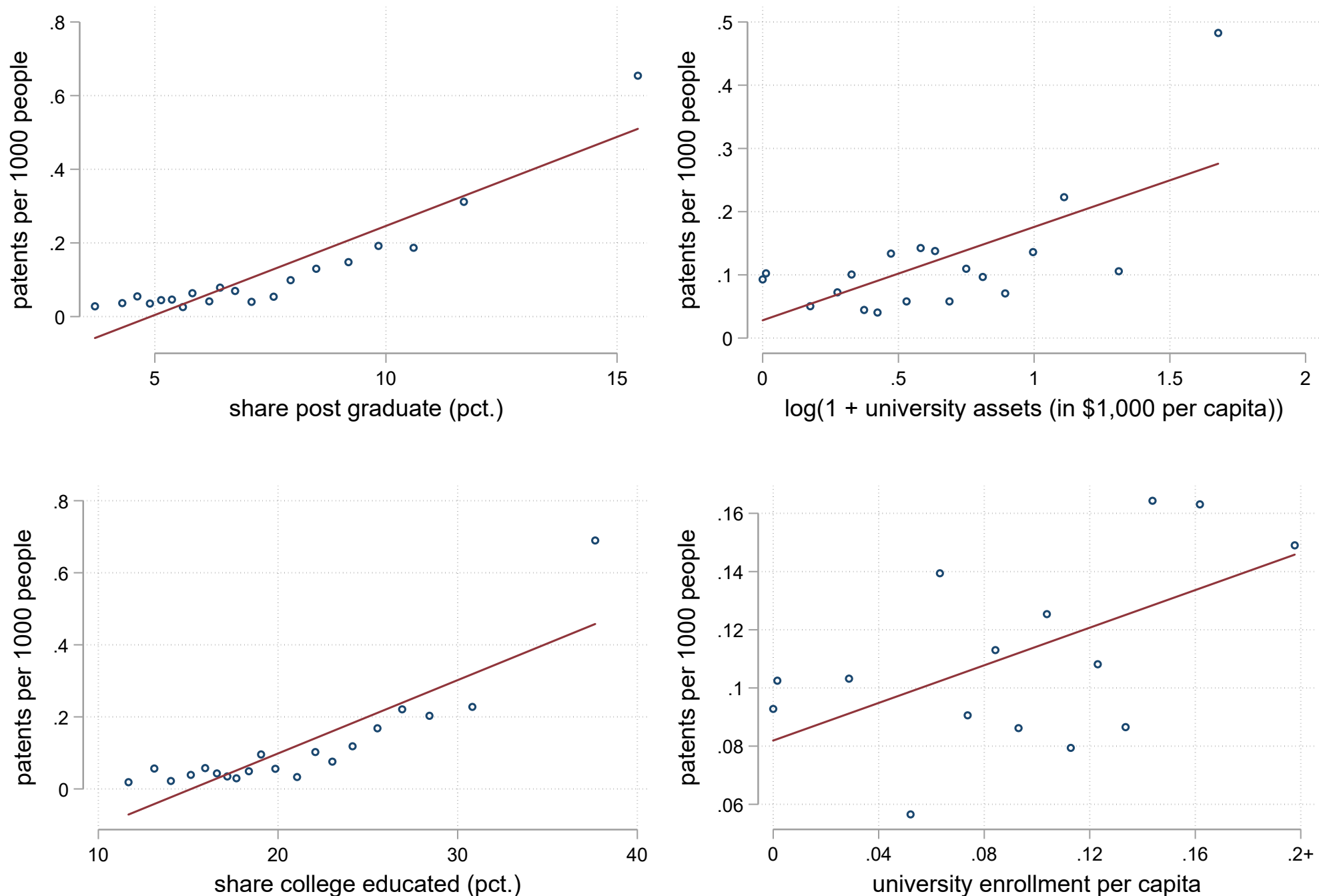
**Appendix Figure 5– Disruptive vs. overall patents, by top CBSAs.**



**Notes:** The figure shows, for disruptive patents (in red) and overall patents (in blue), the normalized share of patenting for the top 20 CBSAs. The normalized share of patents for a CBSA is defined as the share of total patents filed by U.S. inventors in the CBSA (between 1992 and 2016) divided by the share of U.S. population in the CBSA (as of 2015). The figure is sorted by largest to smallest normalized share of disruptive patenting.

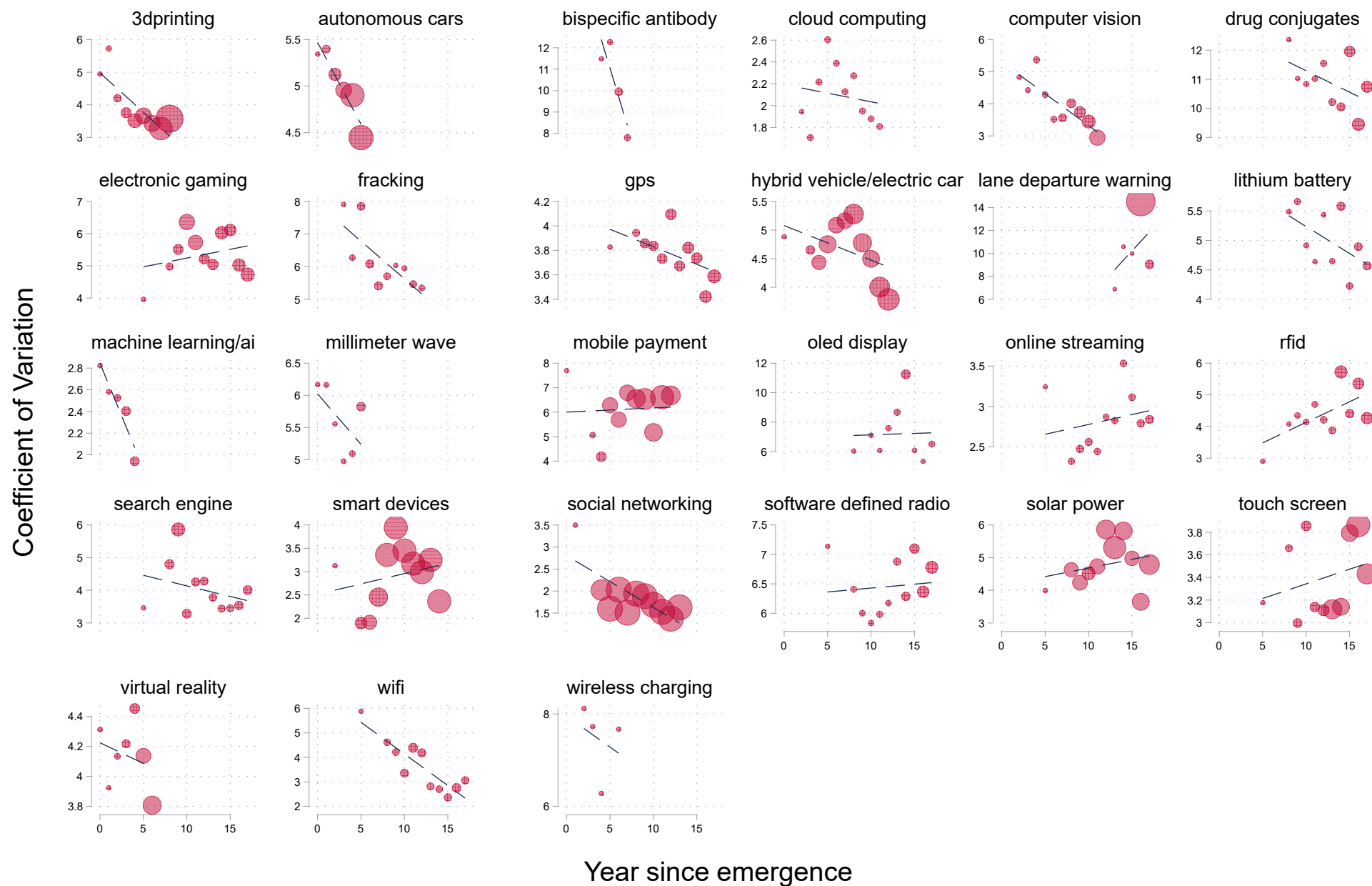


**Appendix Figure 6 - Disruptive innovation vs local skill composition**



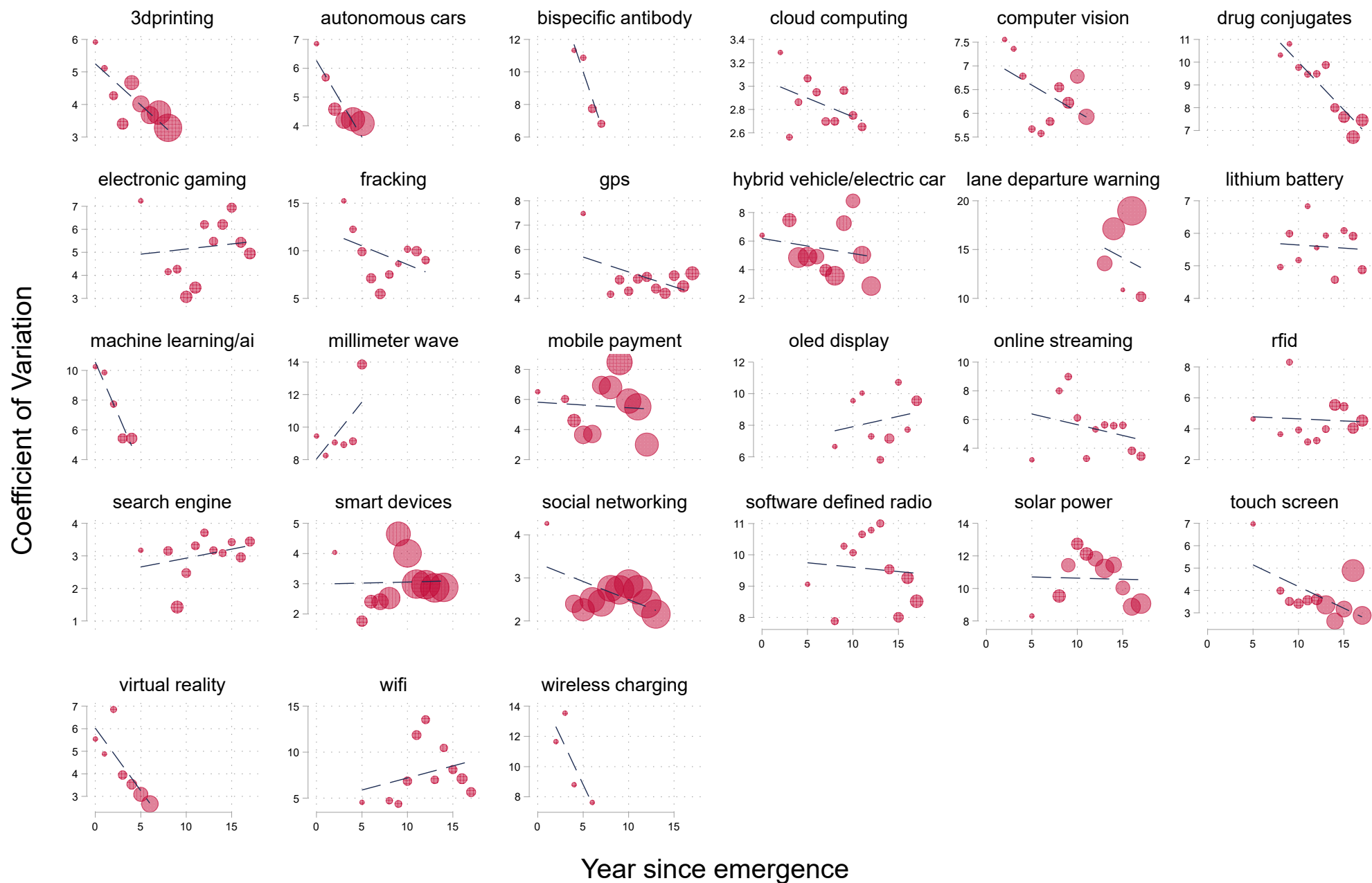
**Notes:** The figures show binned scatter plots of patents associated with a disruptive technology per 1000 people in a CBSA for each of our 29 technologies over measures of skills and university presence in the CBSA. The patents associated with a technology are calculated during the 10 years before the year of emergence of the technology. The fitted lines control for technology fixed effects.

**Appendix Figure 7 – Coefficient of Variation across industries, by year since emergence**



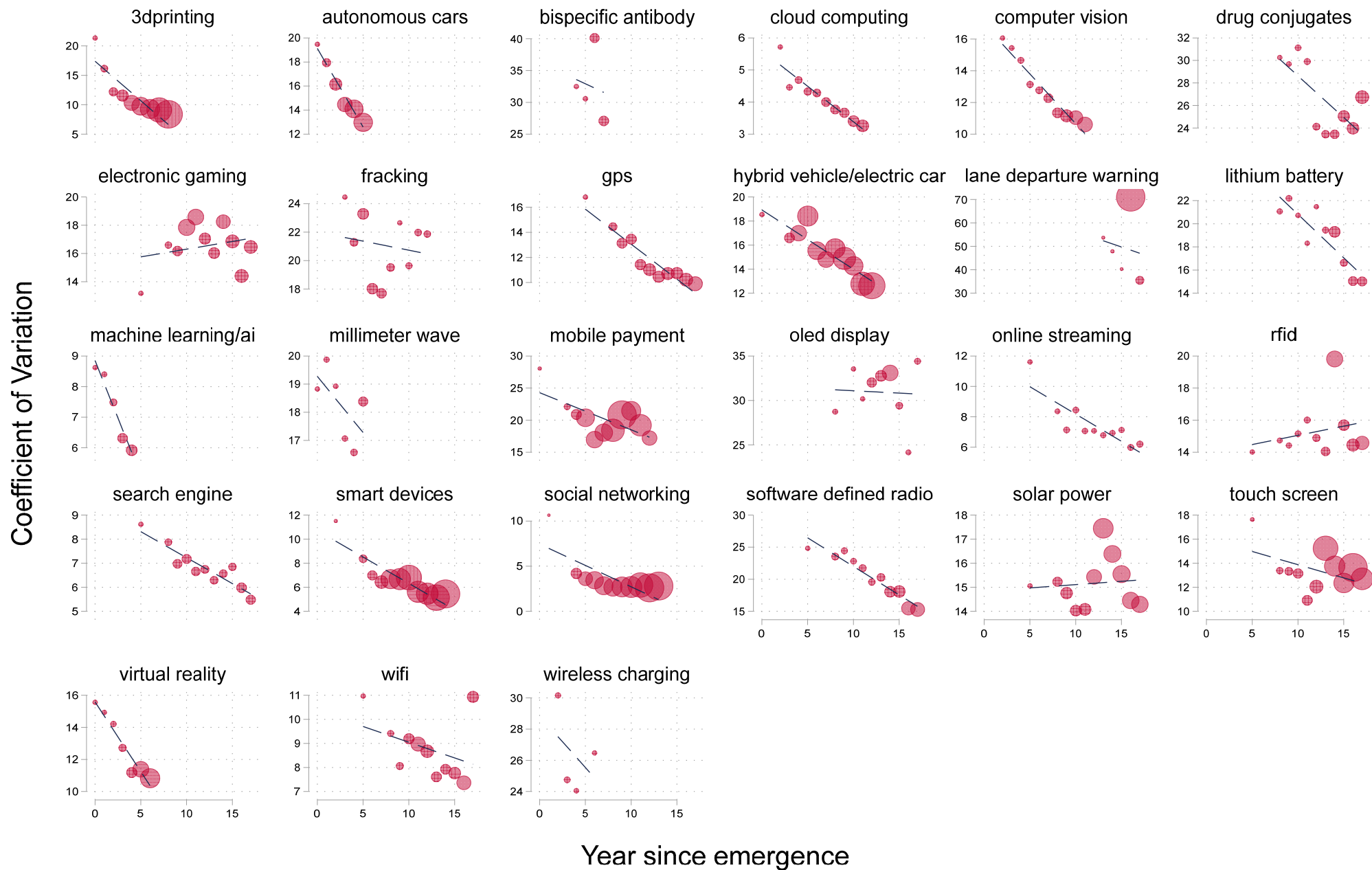
**Notes:** The figure plots the coefficient of variation measured over the normalized share of technology job postings for each of 29 technologies by year from 2007 to 2019 over the years since emergence of the technology.  $Normalized\ share_{i,\tau,t} = \frac{share\ jobs\ exposed_{i,\tau,t}}{share\ jobs\ exposed_{\tau,t}}$ , where  $i$  is an industry.

**Appendix Figure 8 – Coefficient of Variation across occupations, by year since emergence**



**Notes:** The figure plots the coefficient of variation measured using the normalized share of technology job postings for each of 29 technologies by year from 2007 to 2019 over the years since emergence of the technology.  $Normalized\ share_{i,\tau,t} = \frac{share\ jobs\ exposed_{i,\tau,t}}{share\ jobs\ exposed_{\tau,t}}$ , where  $i$  is an occupation.

**Appendix Figure 9 – Coefficient of Variation across firms, by year since emergence**



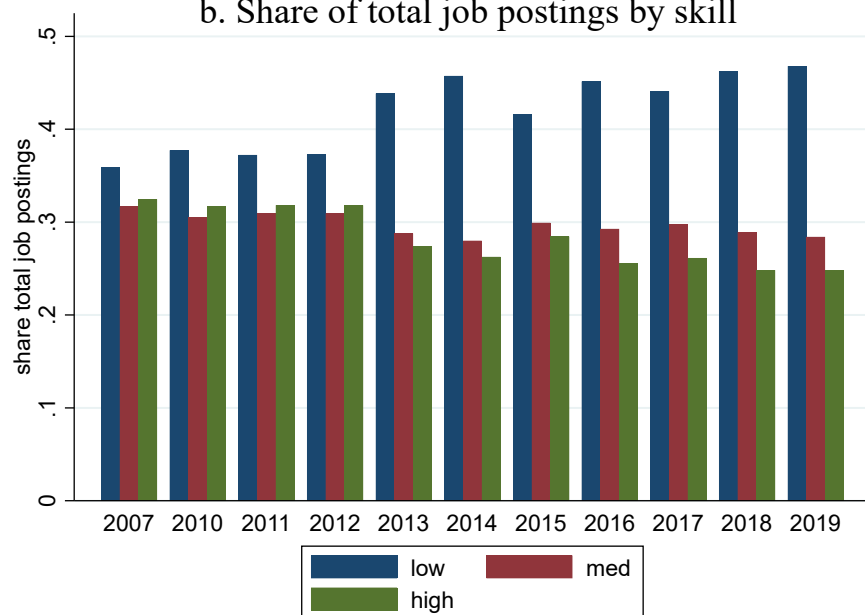
**Notes:** The figure plots the coefficient of variation measured using the normalized share of technology job postings for each of 29 technologies by year from 2007 to 2019 over the years since emergence of the technology.  $Normalized\ share_{i,\tau,t} = \frac{share\ jobs\ exposed_{i,\tau,t}}{share\ jobs\ exposed_{\tau,t}}$ , where  $i$  is a firm.

## Appendix Figure 10 – Overall patterns of Burning Glass job postings

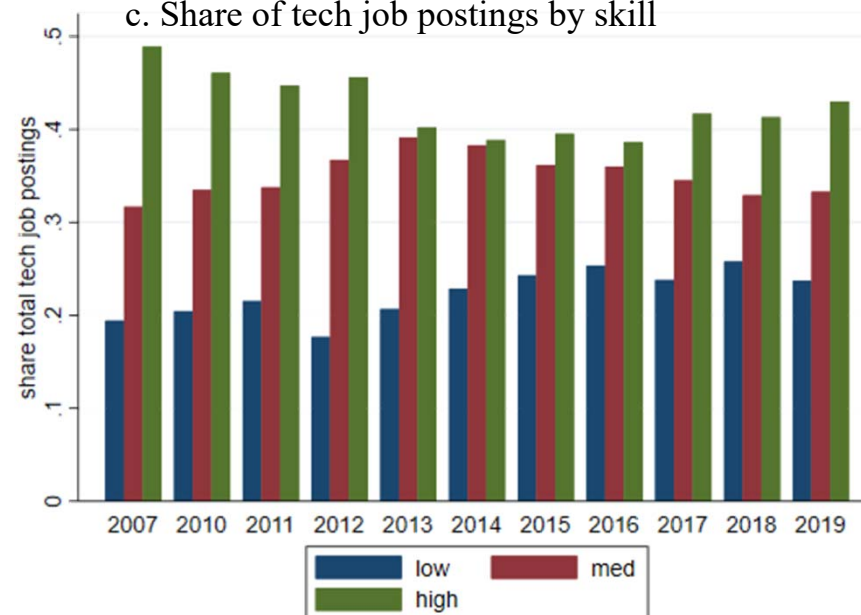
a. BG total job postings (in million)



b. Share of total job postings by skill

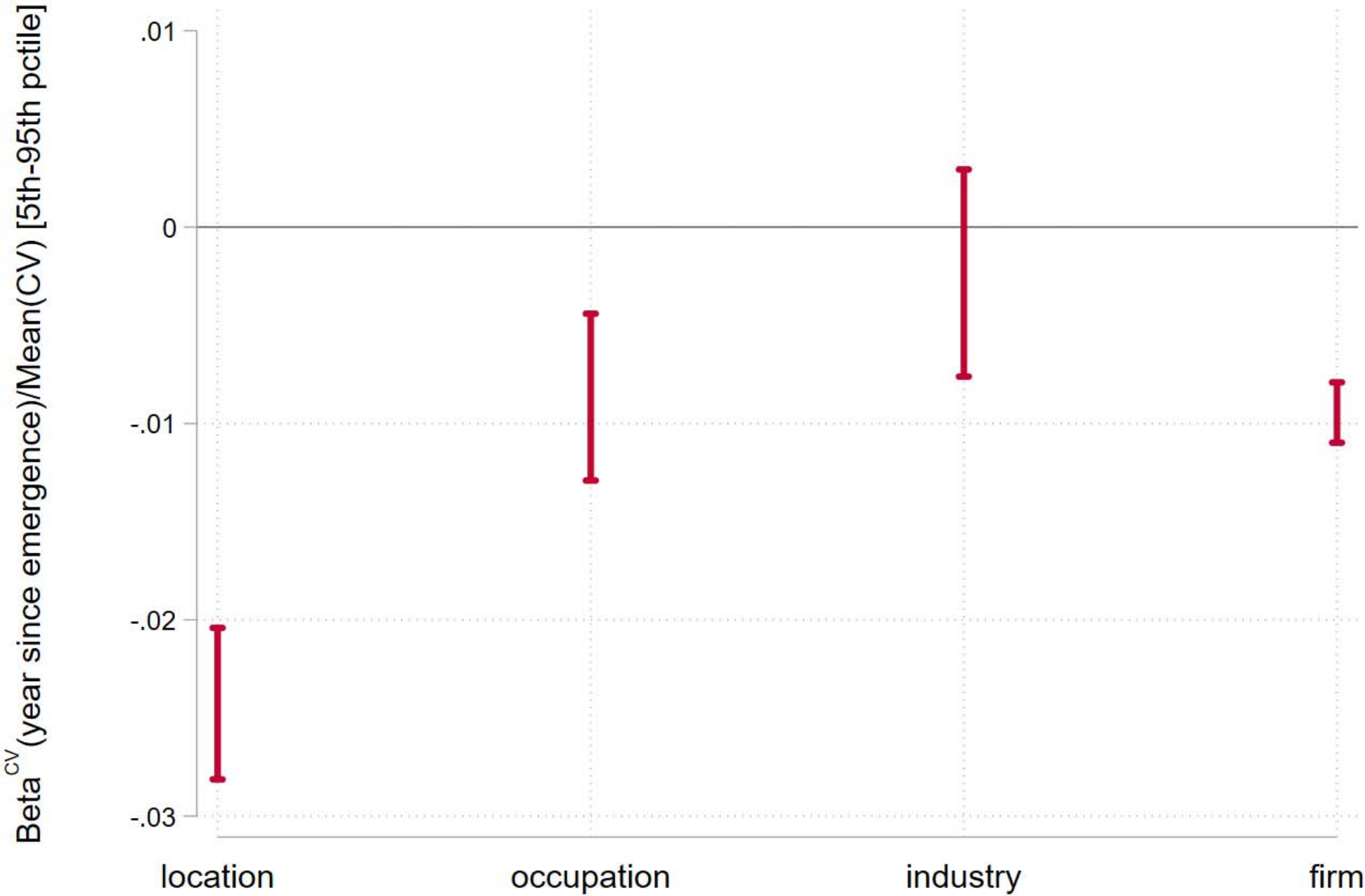


c. Share of tech job postings by skill



**Notes:** In this figure, we show aggregate patterns of Burning Glass (BG) online job postings. Panel a shows the total number of job postings (in millions) by year in BG. Panels b and c show the share of total job postings by skill and the share of total technology job postings (aggregated over 29 selected technologies) by skill, respectively. To calculate skill level for job postings, we aggregate the data over occupations, and then use the share of college-educated workforce from the 2015 American Community Survey to assign them to high-skill occupations (with the share of college-educated people > 60%), medium-skill occupations (with the share of college-educated people > 30% and < 60%), and low-skill occupations (with the share of college-educated people < 30%).

Appendix Figure 11 – Region broadening regressions, jackknife robustness



**Notes:** This figure plots results from a jackknife estimate in regressions of the coefficient of variation of the normalized share of technology job postings calculated across locations, occupations, industries, and firms. For our jackknife estimates, we exclude three technologies at a time and recalculate the degradation in the coefficient of variation. This provides us with 7,308 permutation estimates. In this figure, we plot we plot the 5th and 95th percentile of these jackknife estimates.

## **Data Appendix**

We process four sources of text data, and then combine them with census and university data to conduct our analyses. In this section, we first describe the sources of text data and then the additional auxiliary datasets.

### **1. Sources of text data**

#### **1.1. Earnings conference call transcripts**

From Refinitiv EIKON, we collect the complete set of 321,189 English-language transcripts of earnings conference calls held from 2002 through 2019. Out of these, we drop 5,552 transcripts because we could not reliably match them to a company name in Compustat. We obtain a total of 11,992 firms and 301,294 firm x quarter observations. For our analysis, we aggregate this data up to firm x year level.

#### **1.2. Patents**

We download two separate sets of patent award data for about six million utility patents applied for at the US Patent Office (USPTO) between 1976 and 2015. First, we download full patent text XML files from the USPTO website. Second, we download processed patent variables, such as assignee names, inventor names, location, application year, citations (through 2018), and technology classes from PatentsView.org. We map the FIPS county identifier provided for inventor of each patent to Core-Based Statistical Areas (CBSAs) using a crosswalk provided by the US Census Bureau. For patents with multiple inventor CBSAs, we assign the patent to each CBSA. We also standardize citation counts to control from truncation and time differences: we divide citations for each patent by the average number of citations for patents in the same three-digit CPC technology class and application year.

Furthermore, we map patents to Compustat firms and NAICS three-digit industries in three steps. First, we use the match between patent assignees and public firms in Compustat using the dataset provided by Autor et. al. (2020). Second, since patents themselves do not contain industry codes, we assign them the individual NAICS four-digit industries of the mapped Compustat firms. At the end of this step, we allocate these patents, which represent 45.31% of the original sample, to NAICS three-digit industries.

#### **1.3. Corpus of Historical American English**

Corpus of Historical American English (COHA) is a collection of 116,759 documents published between 1880 and 1970. These include fiction and non-fiction books, and newspaper and magazine articles. As with patents and earnings calls, we decompose these documents into about 400 million unique bigrams.<sup>1</sup> We call these “non-technical” bigrams.

#### **1.4. Burning glass job postings**

From Burning Glass (BG), we obtain about 200 million job postings posted online in the US. Similar to patents, job postings data is provided in two sets. The first contain the full text coded in XML files. We undertake minimal processing of these job postings’ textual data: (1) removing non-letter sections of job postings; (2) removing the top 50 and bottom 50 words from each job posting, as mentioned in Section 2; and (3) as a consequence of step (2), excluding any job posting with less than 100 words. We then perform word counts over the remaining text.

Second, BG codes job postings into occupations (Standard Occupational Classification (SOC) codes), locations (counties), and industries (North American Industrial Classification Codes). BG also extracts an employer name with these job postings. Data Appendix Table 1 provides coverage of these variables in Burning Glass.

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

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

- 1) We search the raw employer string on Bing.com and store the top five search result links. E.g., for the employer string “Tesla Incorporated,” we get <https://www.tesla.com>, <https://en.wikipedia.org/wiki/tesla-inc>, <https://www.britannica.com/topic/tesla->

---

<sup>1</sup> COHA was downloaded from [www.english-corpora.org/coha](http://www.english-corpora.org/coha).



motors, <https://www.bloomberg.com/quote/tsla/us>, <https://www.marketwatch.com/investing/stock/tsla>.

- 2) We group two employer strings under a single identifier if they share at least two out of top five links in common with each other.

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

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

- a. Decompose employer and assignee strings into 5 letter combinations. For example, “Alphabet” is broken into: “alpha,” “lphab,” “phabe,” and “habet.”
- b. Calculate a term frequency, which is the frequency of the five letter combination in the string. In our example of “Alphabet”, each combination uniquely appears in the strings. We calculate an inverse document frequency (idf), which is inverse of the frequency with which the combination appears in all strings of assignees and BG employers.
- c. We combine the term frequency (tf) with an inverse document frequency (idf) to obtain a vector of combinations for each string:

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

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

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

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

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

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

## **2. Processing text to obtain Disruptive Technologies**

As explained in the paper, we use patent text, the COHA, and earnings calls to get to our list of disruptive technical bigrams. The process is described here in detail.

### **2.1. Patents to technical bigrams:**

A typical patent award has five text sections: (1) Title, (2) Abstract, (3) Background, (4) Detailed Descriptions, and (5) Claim. We combine text from all of these sections into one large text string, and then break it down into two-word combinations or bigrams. In this process, we read only those bigrams which are mentioned at least twice in a patent and are not mentioned in “non-technical” bigrams obtained from the COHA.

We then sort these bigrams in terms of their importance to all patents. To do so, we collect the (standardized) number of citations for each patent in our sample and allocate it to each bigram that is mentioned in these patents. As an example: if a patent 123 has bigrams “A” and “B” and is cited twice, then we attribute two citations to bigrams “A” and “B.” Through this process, we cumulatively add up all citations attributed to the “technical” bigrams mentioned in our sample of patents. We end up with a cumulative citations value for 1,509,306 technical bigrams.

In order to focus our attention only on the most important bigrams, we only keep bigrams with at least 1,000 cumulative standardized citations. At the end of this process, we obtain 35,063 bigrams.

### **2.2. Technical bigrams to disruptive technical bigrams using earnings calls:**

We use textual data from firms’ earnings calls to measure disruption in influential technical bigrams. To do so, we first count the mentions of the 35,063 bigrams identified in the last step in our sample of earnings calls transcripts, and keep those that are mentioned in at least 100 earnings calls. Out of these we find 19,897 bigrams which are mentioned at least once in earnings calls, and 2,181 bigrams are mentioned in at least 100 transcripts. Second, for each year between 2002 and 2019 and for each bigram, we calculate the percentage of firms in our earnings calls’ sample that mention a given bigram. Third, for each bigram, we compute the maximum of this percentage across all years and compute the ratio between the percentage in 2002 and the maximum. This provides us with a degree to which bigram mentions have changed in our sample of firms. Finally, we identify a list of 305 bigrams for which this ratio is 0.10 or less, which means that bigram

mentions have increased in earnings calls by at least 10 times since the first year of observation in 2002.

### **2.3. Manual intervention:**

So far, our process of getting to the list of disruptive technical bigrams has been automated. However, to get to the desired list of technologies, we need to take a series of subjective decisions. This process is described as a supervised approach in the paper. We start with our list of 305 bigrams from the last step and process it into disruptive technologies:

### **2.4. Removing non-technology keywords:**

We manually remove bigrams which refer to (1) economic, engineering, and social problems (such as “carbon footprint” or “power outage”), (2) older technologies (“nand flash”), or (3) any bigram that is vague or refers to multiple innovations, such as “flow profile”. We classify bigrams into problems or older technologies by reading Wikipedia pages that mention these bigrams. We classify bigrams as vague by reading their excerpts in earnings calls transcripts and patents text. At the end of this process, we get a list of 105 bigrams.

### **2.5. Grouping bigrams into technologies:**

We group the shortlisted bigrams into technologies by reading Wikipedia pages that mention these bigrams. If a Wikipedia page is not available, then we turn to the Wikipedia page to which we are redirected by search engines. For example, the bigrams “mobile devices,” “smart phones,” and “mobile platform” all refer to “smart devices,” which Wikipedia defines as “an electronic device, generally connected to other devices or networks via different wireless protocols”. This manual grouping of bigrams provides us with a list of 29 technologies.

We tried automating our grouping exercise by clustering bigrams using search engine results and embedding vectors trained on earnings calls transcripts, patents, and Wikipedia pages. However, all of these automated approaches had a false positive rate of about 10-20%.

### **2.6. Extending the list of bigrams:**

Finally, we extend the list of keywords for each technology in two steps. First, we use Wikipedia pages. In particular, the first paragraph of the Wikipedia pages associated to a given technology usually mentions a set of terminology used to refer to this technology. We take this set as is and add it to our existing list. Second, we use bigram embeddings trained on earnings calls transcripts,

so that we capture words employed in similar context to that of earnings calls. Importantly for our context, bigram embeddings are machine learning models that, after trained in a given sample, provide a similarity measure between two bigrams by using the context in which they are mentioned. We obtain a list of the top most similar bigrams to our 29 technologies by adding up similarity scores across all existing bigrams for the technology. For example, top five most similar bigrams (along with their similarity scores) to initial bigrams grouped into the “Smart Devices” technology in step b) are: “iot devices” (0.67), “smart tvs” (0.64), “handheld devices” (0.64), “portable devices” (0.63), and “smartphone tablets” (0.63).

### **2.7. Counting bigrams in Burning Glass and patents:**

Having shortlisted a list of bigrams for each technology in step (3), above, we count these bigrams in more than 200 million BG job postings. We assign an exposure dummy of 1 if the posting mentions a particular technology. This gives us a dataset that contains a job identifier and whether the particular posting is exposed to any of the 29 technologies. As mentioned earlier, BG provides job text separately from job characteristics (such as occupation, location, and employer): the two files are linked via a unique job identifier. Thus, we perform a merge over 200 million job identifiers and then aggregate technology exposure over occupation, location, firm, industry, technology, and time. At the end of this process, we have a dataset with  $i \times \text{technology} \times \text{time}$  dimensions, where  $i$  is one of occupation, location, firm, or industry.

We use a similar process to count our shortlisted bigrams in the million US patents, and then use the corresponding information on occupation, location, firm, and industry of a patent to aggregate upwards. After aggregating, we have a dataset with a total count of patents for  $i \times \text{technology} \times \text{time}$  cells, where  $i$  is one of occupation, location, firm, or industry. Finally, we use these patent counts to identify respective pioneers along each of the four dimensions, as explained in detail in the paper.

## **3. Auxiliary Data**

We combine the above text data with the following sources of data for occupational, geographic, firm, and industry characteristics.

### **3.1. American Community Survey (2015)**

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

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

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

### **3.2. University Data**

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

- 1) Number of research universities in a CBSA: HERD provides details of universities which spend more than \$150,000 in research. We map university zip codes, provided as a part of university addresses, to CBSAs using a crosswalk provided by US Census Bureau. Finally, we count the number of research universities in a CBSA to construct our variable.

- 2) University assets in a CBSA: IPEDS provides details of finances for most post-secondary educational institutions in the US. As with 1) above, we assign these universities to CBSAs and then aggregate their assets over CBSAs.

### 3.3. Weighing scheme to match BG postings to US hiring

In this section, we calculate a weighing scheme to match BG postings at the two-digit SOC-by-year level to US hiring at the same level. To do this, we first calculate US hiring at the occupation level using a combination of US Census' Longitudinal Employer Household Dynamics (LEHD) and U.S. Bureau of Labor Statistics' Occupational Employment and Wage Statistics (OES) databases. We do this in two steps:

- a. From the LEHD, we download the number of new hires which retain full quarter employment, which is defined as “estimated number of workers who started a job that they had not held within the past year and the job turned into a job that lasted at least a full quarter with a given employer.” This is downloaded for NAICS three-digit industries in each year between 2007 and 2019. We aggregate these numbers up at NAICS three-digit industry-by-year level.
- b. From the OES, we download employment data at the three-digit NAICS by two-digit SOC code-by-year level. We calculate the share of employment in each occupation for workings in each three-digit NAICS industry in each year.
- c. We merge the two datasets and calculate hiring at the occupation using the following formula:

$$H_{o,t} = \sum_i \vartheta_{i,o,t} H_{i,t}$$

where  $H_{i,t}$  is hiring for industry  $i$  at time  $t$ ,  $\vartheta_{i,o,t}$  is the share of employment in industry  $i$  at time  $t$  accounted for by occupation  $o$ , and  $H_{o,t}$  is the calculated hiring for occupation  $o$  at time  $t$ . This calculation assumes that hiring in an occupation within an industry is in proportion of its employment in the same industry.

- d. Finally, we reweigh BG job postings to US hiring using the following formula:

$$J^r_{o,t,\tau} = w_{o,t} J_{o,t,\tau}$$

where  $J_{o,t,\tau}$  denotes job postings for occupation  $o$ , technology  $\tau$  at time  $t$  and  $J^r_{o,t,\zeta}$  denotes the reweighted version of that, and;  $w_{o,t} = \frac{H_{o,t}}{J_{o,t}}$  denotes the weights at the SOC 2-digit x time level which are a ratio of US hiring ( $H_{o,t}$ ) and BG job postings ( $J_{o,t}$ ).