

# Reliance on Scientific Ideas in Patenting<sup>1</sup>

Sam Arts<sup>2</sup>

Department of Management, Strategy and Innovation  
Faculty of Economics and Business, KU Leuven

[sam.arts@kuleuven.be](mailto:sam.arts@kuleuven.be)

Nicola Melluso

Department of Business and Management

LUISS University

[nmelluso@luiss.it](mailto:nmelluso@luiss.it)

## *Preliminary and incomplete*

### ABSTRACT

Scientific ideas drive firm innovation, yet measuring firms' reliance on science remains a challenge. Prior research has often relied on patent-to-paper citations, whose primary purpose is to delineate the legal scope of a patent rather than to document all scientific ideas on which it builds. We develop a novel text-based approach that identifies when patents reuse specific ideas from the scientific literature. We validate this method, including using firm R&D lab-level survey data, and show that patents frequently build on prior scientific ideas without citing the original source paper, particularly when the paper is older or less similar in content. Finally, we demonstrate that reliance on scientific ideas positively correlates with the private value of patents, even after controlling for traditional patent-to-paper citations and within the subset of patents without such citations. Together, these results suggest that traditional citation-based measures underestimate the role of science and illustrate the complementary value of our text-based approach and open dataset for studying how scientific knowledge contributes to firm innovation.

---

<sup>1</sup> We gratefully acknowledge insightful feedback and suggestions from Myra Mohnen and Scott Stern. We also thank Michael Roach and Wes Cohen for providing access to the Carnegie Mellon Survey on Industrial R&D data and for their helpful guidance on its use.

<sup>2</sup> Corresponding author: Sam Arts, Department of Management, Strategy and Innovation, Faculty of Economics and Business, KU Leuven, email: [sam.arts@kuleuven.be](mailto:sam.arts@kuleuven.be), mailing address: Hendrik Conscienceplein 8 - box 15520, 2000 Antwerp, Belgium, phone: +32 3 201 18 36.

## 1. Introduction

Scientific ideas fuel technological progress and economic growth in industry (Bush, 1945; Mokyr, 2011). For instance, Stanley Cohen’s scientific discovery of epidermal growth factor, a cell-growth protein that earned him the Nobel Prize, was foundational for the development of cancer therapies and blockbuster drugs in the pharmaceutical industry.<sup>3</sup> Measuring and studying how firms draw on scientific knowledge to develop new technologies is therefore a central concern in strategy and innovation research (Cohen & Levinthal, 1990; Fleming & Sorenson, 2004). To date, most empirical work has relied on two primary data sources: firm surveys (e.g., Mansfield, 1991; Cohen et al., 2002) and patent citations to scientific publications (e.g., Narin et al., 1997; Krieger et al., 2024). While surveys provide direct and valuable insights, they are typically constrained by small, cross-sectional samples. More recently, patent citations to scientific papers have emerged as the standard large-scale proxy for science–technology linkages (Marx & Fuegi, 2020). Yet, like any data source, patent-to-paper citations have limitations, offering only a partial view of how firms draw on scientific knowledge to drive innovation performance and build competitive advantage (Roach & Cohen, 2013).

These limitations primarily stem from the fact that the main purpose of patent citations is to delineate the legal scope of an invention, not to comprehensively document all sources of prior knowledge—such as scientific ideas—that contributed to the development of the patented technology (Jaffe et al., 1993). As a result, patent-to-paper citations underestimate the role of science in firm innovation (Type II error). Inventor surveys show that 30% to 65% of patents whose inventors report relying on scientific research nonetheless contain no citations to scientific papers (Tijssen, 2002; Callaert et al., 2014), and firm-level surveys suggest that patent-to-paper citations capture only about half of science’s actual contribution to firms’ innovation performance (Roach & Cohen, 2013). At the same time, citations are also prone to Type I error: inventor surveys show that roughly 50–70% of the scientific publications cited in patents were not important for the development of the patented technology (Callaert et al., 2014; Bryan et al., 2024).<sup>4</sup> Finally, patent-to-paper citations rest on

---

<sup>3</sup> The term *epidermal growth factor* first appeared in a 1967 paper co-authored by Stanley Cohen and has since been reused in more than 23,000 U.S. patents, illustrating the far-reaching technological influence of this scientific discovery.

<sup>4</sup> This limitation of citations as indicators of knowledge flows extends beyond patent-to-paper citations. In the context of patent-to-patent citations, a survey of about 380 U.S. inventors found that roughly half of the citations did not reflect any actual knowledge flow or influence on the citing patent’s development (Jaffe, Trajtenberg & Fogarty, 2000). Likewise, in the context of paper-to-paper citations, a survey of 9,000 academics across 15 scientific fields found that more than half of the papers they cited were judged to have had only minor or very minor intellectual influence (Teplitskiy et al., 2022).

a linear, one-way model of innovation in which science is treated solely as an upstream input to technology. This view neglects the dynamic, two-way relationship in which science and technology may advance together through iterative feedback loops (Rosenberg, 1976; Stokes, 1997; Barbosu & Teodoridis, 2025). As Paula Stephan (2012: 84) observes, “scientific research leads to advances in technology, but it is new technology that often brings about advances in science.”

In this paper, we develop and validate a new method to measure a patent’s reliance on scientific ideas and make the underlying data publicly available. Rather than relying on citations, our approach analyzes the full text of all U.S. patents to identify the scientific ideas embedded in each invention. We build on Arts et al. (2025a), who processed the complete historical corpus of scientific publications across all domains to identify novel scientific ideas—captured as newly introduced noun phrases such as *electroporation*, *monoclonal antibody*, *scanning electron microscope*—that first appeared in the scientific literature.<sup>5</sup> By processing the text of both scientific publications and all U.S. patents, we identify, for each patent, the reuse of scientific ideas originally pioneered in the academic literature, as well as ideas that first appeared in patents before diffusing into scientific discourse.

To validate our approach, we first show that patents with clear ties to science—such as those assigned to universities or invented by scientists—score higher on our metrics of reliance on scientific ideas, even in a sample of patents without patent-to-paper citations, which accounts for 78% of all patents. Second, we look for traces of evidence that knowledge was indeed transferred from the source paper introducing a scientific idea to the patent reusing that idea. Patents that reuse a given scientific idea are more likely than matched control patents to (indirectly) cite the scientific paper in which the idea first appeared, and their inventors are geographically and socially closer to the authors of the source paper. Finally, we use firm-level survey data from the Carnegie Mellon Survey on Industrial R&D to show that R&D labs whose managers report greater reliance on scientific research from universities and public research institutions also generated patents that rely more heavily on scientific ideas according to our new metrics, including within the subsample of labs whose patents contain no scientific references, a group that accounts for 54% of all R&D labs.

---

<sup>5</sup> Arts et al. (2025) validate this text-based approach in multiple ways and show that new noun phrases consistently outperform all other novelty metrics—including embedding-based measures that compute the semantic distance between a paper’s text and all prior scientific work—in identifying new scientific ideas and assessing a paper’s novelty at the time of publication. They demonstrate, for instance, that the noun-phrases identify the first appearance of major scientific ideas in the literature and that, in a subset of cases, the detected ideas align with Nobel Prize-winning contributions as independently documented in the official Nobel Prize materials.

To evaluate how accurately patent-to-paper citations capture a patent’s actual reliance on scientific ideas, we use our text-based approach to construct several datasets in which there is strong evidence that a patent builds on a specific scientific idea introduced in a particular paper. These include cases where inventors had previously cited the paper in patents or publications, or had co-authored the paper themselves. Across all settings, citation rates remain relatively low—ranging from 26% to 56%—suggesting that patents often rely on scientific ideas without citing their original source. In line with the legal function of patent references to delineate the scope of protection, we find that patents are more likely to cite the specific scientific papers whose ideas they reuse when those papers are more recent or when their content is more closely aligned with the patent.

Finally, to demonstrate the value of our new approach for strategy and innovation research, we revisit the long-standing question of whether firm inventions that build on science have greater private value, measured by stock market reactions or received citations (Fleming & Sorenson, 2004; Krieger et al., 2024). We find that patents drawing on scientific ideas yield higher private value, a result that holds both after controlling for traditional patent-to-paper citations and when restricting the sample to patents without such citations. This finding underscores how our text-based approach and open dataset provide a powerful complement to patent-to-paper citations in understanding how science fuels firm innovation.

## **2. Identifying patents’ reliance on scientific ideas**

### ***2.1 Data and sample***

Our analysis combines data from multiple sources. For patent data, we use PatentsView, from which we extract the full text—including the title, abstract, claims, summary, detailed description, and drawing descriptions—of all U.S. patents granted between January 1976 and June 2023 (n=7,699,586). For scientific publications, we use the January 2024 snapshot of OpenAlex, as processed and cleaned by Arts et al. (2025a), which includes 75,295,921 papers and conference proceedings published between 1901 and 2023.<sup>6</sup> For brevity, we refer to both as papers throughout the manuscript. We supplement these sources with additional data: assignee data from DISCERN 2.0 (Arora et al. 2024) and PATSTAT, patent-to-paper

---

<sup>6</sup> Arts et al. (2025a) retain papers with an English title and abstract, complete publisher information, a valid impact factor for the publication year, non-empty titles, at least one listed author, and no duplicate titles or abstracts.

citations from Marx & Fuegi (2020), inventor–author linkages from Scharfmann et al. (2024), and patent commercial value estimates from Kogan et al. (2017).

To identify scientific ideas, we begin with the processed text of scientific publications from Arts et al. (2025a), which yields 27,079,343 distinct noun phrases.<sup>7</sup> Although they use a baseline dictionary and apply preprocessing to remove boilerplate language and generic terms unlikely to represent meaningful scientific ideas, their method remains inherently imperfect. To further refine the identification of genuine scientific ideas, we use an LLM to assign each noun phrase a confidence score from 1 to 3, reflecting the likelihood that it captures a true scientific idea. In the remainder of the paper, we exclude phrases with a score of 1 (8.9% of the sample). These phrases are rated as unlikely to represent a scientific idea, and validation exercises confirm that excluding them reduces false positives without reducing coverage of true ideas (see Appendix B for details).<sup>8</sup>

We then process the full text of all patents using the same preprocessing and extraction procedures applied to the scientific corpus. Of the noun phrases identified in scientific publications, 29.32% also appear in patents, while 70.68% appear only in scientific papers. Among noun phrases that appear in both papers and patents, 57.13% (n = 4,535,841) first appear in a scientific paper, and 42.87% (n = 3,403,386) first appear in a patent, as determined by the publication date of the paper and the filing date of the patent. As such, a notable share of ideas are first disclosed in patents and only subsequently appear in the scientific literature. Among all phrases first introduced in patents, 23.8% (n = 621,807) subsequently appear in a scientific paper for the first time after the patent’s filing date but before it is granted. Only 1.7% first appear in a paper co-authored by at least one inventor of the pioneering patent, suggesting that only a very small minority are likely to represent patent–paper pairs.

To give one example, *interleukin-7* was first introduced in a patent (US 4,965,195) filed by Immunex (which was later acquired by Amgen for 16 billion USD) in 1988 describing how to clone and produce this immune-system protein, which proved essential for understanding T-cell and lymphocyte development. The patent gave Immunex strong control

---

<sup>7</sup> A noun phrase consists of a single noun or a sequence of words centered on a noun (the head), along with modifiers that refine its meaning. It is widely used in computational text analysis to capture distinct concepts in natural language, making it a practical unit for identifying scientific ideas or concepts in the text of scientific papers.

<sup>8</sup> We provide open access to the full set of noun phrases along with their associated confidence scores.

over interleukin-7-based therapeutics. The discovery has since influenced both scientific immunology research and technological development in clinical therapies for immune reconstitution, as reflected in the reuse of *interleukin-7* in 997 subsequent scientific papers and 1,536 patents. In line with prior work examining how patents and firms rely on scientific ideas, the remainder of the paper focuses on ideas first introduced in the scientific literature and later reused in technology development (Iaria et al., 2018).<sup>9</sup> However, unlike patent-to-paper citations—which are inherently unidirectional—our data and code also enable analysis of the reverse pathway, where new technologies drive scientific progress (e.g., Rosenberg, 1976; Barbosu & Teodoridis, 2025).

Table A.1 in Appendix lists the top ten noun phrases that first appeared in scientific literature and were later most frequently reused in patents, grouped by the decade in which they were first adopted in patents. Early decades are dominated by materials and chemical sciences (e.g., *polypropylene*, *polyethylene glycol*), later decades by computing and electronics (e.g., *liquid crystal display*, *field-programmable gate array*), and the most recent decades by life sciences (e.g., *imatinib*, *cas9*) and advanced information sciences such as artificial intelligence (e.g., *generative adversarial network*, *smart contract*). This progression reflects the historical shift in scientific discovery—and the resulting technological opportunity—from chemistry and materials toward computing, life sciences, and information sciences, including artificial intelligence.

## 2.2 Descriptive statistics

We measure a focal patent’s reliance on scientific ideas or concepts by identifying noun phrases that were first introduced in scientific papers and later reused in the focal patent’s full text. Based on this, we construct three measures: *Science Phrases* as the total number of unique scientific noun phrases in the patent, *Science Phrases Share* as their proportion relative to all unique noun phrases in the patent, and *Science Phrases Bin* which is a binary indicator equal to one if the patent reuses at least one noun phrase first introduced in the

---

<sup>9</sup> A natural alternative is to use embedding-based similarity measures that compute the semantic proximity between a patent’s text and prior scientific publications (Ghosh et al., 2024). However, as shown by Arts et al. (2025a), such methods perform worse in identifying new scientific ideas and in measuring the novelty of a paper at the time of publication. Moreover, unlike patent-to-paper citations or the reuse of scientific noun phrases in patent text, embedding-based approaches cannot reveal which specific scientific ideas a patent actually builds on. Semantic proximity is also not directional: a patent may be close in embedding space to many papers in the same field, even though only one (or none) of them introduced the underlying scientific idea that the patent reuses. By contrast, the reuse of specific scientific noun phrases—each traced to the first paper in which it appeared—directly identifies both the content and the origin of the scientific idea incorporated into the patent.

scientific literature.<sup>10</sup> For each focal patent, we also calculate *Science References* as the number of distinct scientific papers cited by the focal patent on the front page or in the body text of the patent (Marx and Fuegi, 2020; Bryan et al., 2020), and *Science References Bin* which is a binary indicator equal to one if the patent cites at least one scientific paper.<sup>11</sup>

*'Insert Table 1'*

Table 1 reports summary statistics at the patent level. The full text of the average (median) patent contains 255 (179) distinct noun phrases, of which 32 (9) were first introduced in the scientific literature before being reused in the focal patent. Accordingly, *Science Phrases Share* for the average (median) patent is 0.081 (0.053), and 96% of all patents reuse at least one scientific noun phrase. By comparison, patents cite only 3 scientific papers on average (median = 0), and only 22% of patents cite at least one scientific paper. The correlation between *Science Phrases* and *Science References* is 0.58, suggesting that while the two measures are positively related, they also capture distinct dimensions of a patent's reliance on science. As shown in Figure A.1 in Appendix, even patents without any science references, representing 78% of the full population, exhibit substantial variation in *Science Phrases*, suggesting that *Science References* may underestimate a patent's reliance on scientific ideas. In fact, 95% of patents without any scientific references reuse at least one idea (noun phrase) from the scientific literature, and these patents reuse an average of 14 scientific ideas (median = 10).

Notably, the average (median) age of reused scientific ideas—measured as the number of years between the publication of the paper in which the phrase first appeared and the patent's filing year—is 46 (46) years. This is substantially older than the average (median) age of cited scientific papers, measured as the number of years between the publication of the cited paper and the patent's filing year, which is only 10 (9) years. Thus, patents tend to cite

---

<sup>10</sup> One might argue that longer patents naturally contain more *Science Phrases*, potentially biasing the measure. However, our analyses control for patent text length, and results remain robust when using *Science Phrases Share*, which corrects for length by construction. A related concern is that multiple distinct noun phrases in a patent may capture the same underlying scientific concept. Here too, *Science Phrases Share* mitigates this issue by construction.

<sup>11</sup> To measure the number of science references made by a patent, we rely on the same population of scientific papers from OpenAlex that we use to construct our measure of patents' reliance on scientific ideas. Because OpenAlex is very broad and includes working papers and obscure journals, we follow Arts et al. (2025) and restrict the sample to journal articles and conference proceedings that (i) have an English title and abstract, (ii) include complete publisher information, (iii) are linked to a valid impact factor in the publication year, (iv) contain non-empty titles, (v) list at least one author, and (vi) are not duplicate records. Importantly, the correlation between our restricted measure of science references and an unrestricted count using the full OpenAlex sample (as in Marx and Fuegi, 2020) is 0.998, and all results reported in this paper hold under both approaches.

more recent scientific papers, even though the underlying inventions often also rely on older, foundational scientific ideas without citing the original source publications.

Table A.2 lists, for selected industries, the top five U.S. public firms whose technology development most heavily relies on science, as measured by the average number of distinct scientific phrases reused per patent, along with each firm's five most frequently reused scientific phrases. In the Chemicals and Allied Products industry (including Biotech), Genentech's most reused scientific phrases—*polypeptide*, *amino acid sequence*, and *plasmid*—reflect its deep roots in molecular biology research underlying recombinant protein and antibody therapies, while Amgen's *monoclonal antibody* and *polyethylene glycol* capture its reliance on immunology and polymer chemistry for engineering and stabilizing biologic drugs. Ionis Pharmaceuticals' most reused phrases—*oligonucleotide*, *phosphorothioate*, and *nucleoside*—derive from nucleic-acid chemistry and form the scientific basis for its RNA-targeted therapeutics. In the Measuring, Analyzing, and Controlling Instruments industry, Hologic's reuse of *nucleic acid amplification* and *oligonucleotide* reflects its use of genetic-testing research to develop molecular diagnostics, while Thermo Fisher and Waters build on analytical chemistry—especially *mass spectrometry* and *liquid chromatography*—to design instruments for chemical and biochemical analysis. Bio-Rad's heavy reuse of *peptide* and *oligonucleotide* likewise reflects its reliance on protein chemistry and nucleic-acid methods to produce reagents, assay kits, and life-science instrumentation. In the Petroleum Refining industry, ExxonMobil's frequently reused phrases—*copolymer*, *polyolefin*, and *molecular sieve*—stem from polymer science and catalysis research that support its development of high-performance plastics, advanced materials, and catalytic processes for fuel refining.

### **3. Validation of text-based metrics**

We validate our new text-based method using different approaches. First, we assess face validity by testing whether patents with clear links to science—such as those that cite scientific publications, are assigned to universities, or list scientist-inventors—exhibit greater reliance on scientific ideas according to our measures. Second, we validate the method's ability to capture inventors' and patents' reliance on specific scientific ideas by examining whether patents that reuse a given scientific idea are more likely than matched-control patents to (in)directly cite the scientific paper in which that idea first appeared, and whether the inventors on these patents are geographically or socially closer to the authors of the source paper. Third, we use firm-level survey data from the Carnegie Mellon Survey on Industrial

R&D to test whether R&D labs whose managers report relying more heavily on public research also produce patents that rely more heavily on scientific ideas.

### **3.1 Face validity**

As a first validation exercise, we compare groups of patents that should, on average, rely more heavily on scientific knowledge—namely, those that cite at least one scientific publication, are assigned to a university, or list a scientist among their inventors—to the rest of the patent population. Table 2 confirms this pattern. Patents in these groups are more likely to contain at least one scientific idea, reuse a larger number of scientific ideas, and exhibit a higher share of scientific content. For example, patents with science references reuse on average 91 *Science Phrases* and have a 15% *Science Phrases Share*, compared to 14 and 6% for patents without science references. University patents show a similar contrast (103 versus 29 *Science Phrases* and 16% versus 8% *Science Phrases Share* on average) as do patents with scientist-inventors (49 versus 15 *Science Phrases* and 10% versus 6% *Science Phrases Share* on average). While *Science Phrases* and *Science References* perform similarly in distinguishing university from non-university patents, as indicated by comparable Mann–Whitney test statistics in Table 2, *Science Phrases* substantially outperform *Science References* in distinguishing patents with and without scientist-inventors. Importantly, as illustrated in Table A.3, even within the subset of patents without any science references, university patents and those with scientist-inventors still exhibit significantly higher *Science Phrases* and *Science Phrases Share*. This indicates that our text-based measure uncovers meaningful residual variation in patents’ reliance on science that patent-to-paper citations fail to detect.

*‘Insert Table 2’*

### **3.2 Validation of patents’ reliance on scientific ideas: citation, geographic, and social proximity between inventors and authors**

As a second validation, we test whether patents that reuse a scientific idea exhibit stronger evidence of knowledge transfer from the source paper to the patent. We begin by drawing a random sample of 100,000 scientific ideas (noun phrases) first introduced in scientific papers and subsequently reused in at least one patent. For each idea, we identify all patents that reuse it, yielding 7,139,906 reusing patent–paper dyads. For each of these, we construct a matched-control patent–paper dyad, linking the same source paper to a patent that does not reuse the

scientific idea but is filed and granted in the same years and assigned to the same technology subclass as the reusing patent.

*‘Insert Table 3’*

As reported in Table 3, patents reusing a scientific idea are 67 times more likely than matched-control patents to directly cite the scientific paper in which that idea first appeared. In addition to direct citations, patents may also reference the same scientific idea indirectly by citing another patent or paper that, in turn, cites the original source paper (Ahmadpoor and Jones, 2017). We find that patents reusing a scientific idea are 20 times more likely to cite a patent that cites the source paper and 17 times more likely to cite a paper that cites the source paper. These patterns confirm that patents drawing on specific scientific ideas are far more likely to reference the source publication—directly or indirectly—supporting the validity of our text-based approach as a measure of patents’ reliance on specific scientific ideas. Nevertheless, only a small fraction of reusing patents cites the corresponding source paper, consistent with Arts et al. (2025a), who show that only a minority of papers drawing on scientific ideas include a citation to the source papers.<sup>12</sup>

However, (in)direct patent-to-paper citations represent only one proxy for knowledge flows from science to technology. Prior research shows that the transfer of scientific ideas into technological invention is both geographically and socially constrained (Arts et al., 2025b; Balsmeier et al., 2025). To further validate our text-based approach, we compute geographic and social proximity for each of the 7,139,906 reusing patent–paper dyads and for each of the 7,139,906 matched-control patent–paper dyads. If our text-based approach genuinely captures inventors’ and patents’ reliance on scientific ideas—rather than superficial textual overlap introduced during the patenting process—then patents reusing a scientific idea should originate from inventors who are geographically or socially closer to the authors of the paper in which that idea first appeared (Verluisse et al., 2025).

As shown in Panel A of Table 3, inventors on patents reusing scientific ideas are located, on average, 303 miles closer to the authors of the corresponding source paper than inventors on matched-control patents, a difference that is both statistically and economically

---

<sup>12</sup> Interestingly, Figures A.2 and A.3 reveal substantial heterogeneity in citation rates across technologies and industries. Citation rates are highest in chemistry, human necessities, and physics, and lowest in fields such as mechanical engineering. At the industry level, industries such as drugs, laboratory instruments, and surgical and medical devices stand out with the highest citation rates, while industries such as motor vehicles and equipment exhibit much lower citation rates.

significant. A similar pattern emerges for social proximity. Using a collaboration network that connects all inventors on U.S. patents and all authors in OpenAlex through prior co-authorships and co-inventorships (Scharfmann et al., 2024), we measure the shortest interpersonal path between any inventor on a patent and any author of the corresponding paper at the time of patent filing. Inventors on reusing patents are systematically closer to the authors of the source paper than those on matched-control patents: reusing patents are 57 times more likely to include an inventor who is also an author of the source paper, 9 times more likely to involve a direct prior collaboration between an inventor and an author, and 2.5 times more likely to have an indirect connection through a shared collaborator. Interestingly, Panel B of Table 3 shows that geographic and social proximity peak in the first years following a scientific discovery. For example, within the first five years after publication, 12 percent of the patents building on a new scientific idea originate from inventors within two degrees of separation from the authors whereas this share declines to under 1 percent over the full observation window. This pattern suggests that new scientific ideas are initially absorbed locally—by geographically and socially proximate inventors—before diffusing more broadly.

Together, these results support the validity of our text-based proxy in capturing genuine inventor- and patent-level reliance on scientific ideas. Patents reusing scientific ideas are not only more likely to (in)directly cite the papers in which those ideas originated but also disproportionately involve inventors who are geographically and socially closer to the scientists behind them.

### ***3.3 Validation using R&D labs' self-reported reliance on science***

A core advantage of firm-level survey data is that it provides a direct measure of how much a firm relies on scientific research when developing new technology (e.g., Mansfield, 1991; Cohen et al., 2002). Although survey data typically cover only a small cross-sectional sample of firms, it offers a unique opportunity to directly validate alternative large-scale measures of firms' reliance on science (Roach and Cohen, 2013; Bryan et al., 2020). Consistent with this line of work, we use the 1994 Carnegie Mellon Survey on Industrial R&D, linked to patent data, to validate our text-based approach for measuring patents' and firms' reliance on science.

The survey asks R&D lab managers to report the share of their unit's R&D projects that used research findings from universities or government labs during the preceding three years (1991–1993). Responses are given on a five-point ordinal scale (<10%, 10–40%, 41–

60%, 61–90%, >90%), and we follow prior work in treating this variable as our dependent variable. Our sample includes all surveyed R&D labs that answered this question and that filed at least one patent during the same 1991–1993 window. For each R&D lab, we retrieve all U.S. granted patents filed between 1991 and 1993 and compute our explanatory variables at the lab level (Roach and Cohen, 2013). *Science References Avg* is the average logged number of scientific references per patent, and *Science Phrases Avg* is the average logged number of scientific noun phrases per patent. All models control for average patent text length and the logged number of patents filed by the lab during 1991–1993. The sample includes 646 R&D labs from 35 different industries, each linked on average to 32 U.S. patents.

*‘Insert Table 4’*

Table 4 reports the results of the ordered logit regressions, estimated both without (Columns 1–4) and with industry fixed effects (Columns 5–8), and, in Columns (4) and (8), for the subsample of R&D labs whose patents contain no patent-to-paper citations, representing 54% of the surveyed R&D labs. Across all specifications, *Science Phrases Avg* is significantly and positively associated with survey-reported reliance on public research.<sup>13</sup> These associations persist even when including *Science References Avg* (Columns 3 and 7) and for the subsample of R&D labs for which none of the patents filed during the prior three years contained any scientific references. Although *Science Phrases Avg* has lower explanatory power (Pseudo R<sup>2</sup>) than *Science References Avg* in the models without industry fixed effects, their explanatory power is nearly identical once industry fixed effects are included. Table A.4 in Appendix shows that our results are robust when using the “open science” factor from Roach and Cohen (2013) as an outcome variable, which is constructed from survey responses in which R&D managers rate, on a four-point Likert scale, the importance of publications, conferences, and informal communication as channels through which public research contributes to their R&D. The three items load onto a single “open science” factor, yielding a continuous measure of a lab’s reliance on open-science channels. Together, these results validate that our text-based measure genuinely captures firm R&D labs’ reliance on scientific ideas and that it explains meaningful lab-level heterogeneity missed by traditional patent-to-paper citations.

---

<sup>13</sup> The only exception is that the coefficient on *Science Phrases Avg* is marginally insignificant ( $p = 0.101$ ) in Column 8, which reflects the subsample of R&D labs without any patent-to-paper citations and includes industry fixed effects, likely due to the relatively small number of observations ( $N = 349$ ). Estimating the same model without industry fixed effects (Column 4) yields a significant coefficient with  $p = 0.047$ .

## 4. Validation of Patent-to-Paper Citations

### 4.1 Do patent-to-paper citations accurately capture reliance on scientific ideas?

To assess how well patent-to-paper citations capture a patent's actual reliance on scientific ideas, we use our new text-based approach to construct three datasets at the paper–patent level in which there is strong indication that the focal patent draws on a specific scientific idea introduced in the corresponding paper. We then examine how often the patent cites that source paper.

First, following Lampe (2012), we identify from the full population of U.S. patents all cases in which a patent  $p$ , invented by inventors  $i$ , builds on a scientific idea (i.e., reuses a noun phrase) that was first introduced in a paper previously cited by any earlier patent co-invented by at least one of the inventors  $i$ . This results in 764,835 reusing patent–paper dyads where the focal patent reuses a scientific idea from the focal paper that at least one of the inventors had cited in a prior patent and thus is arguably aware of. Interestingly, only 56% of these cases include an explicit citation to the original paper.

Second, we conduct a parallel analysis based on citations in inventors' prior scientific publications, rather than in their prior patents. Unlike patent-to-paper citations, which may be added by patent attorneys or examiners, citations in scientific papers are typically included by the authors themselves and thus more directly reflect their knowledge of the cited work. We identify from the full population of U.S. patents all cases in which a patent  $p$ , invented by inventors  $i$ , reuses a scientific idea that was first introduced in a paper previously cited by any earlier scientific paper co-authored by one of the inventors  $i$ . This yields 343,080 reusing patent–paper dyads in which the focal patent builds on a scientific idea introduced in the focal paper previously cited by one of its inventors in their earlier scientific work. Yet, only 26% of these observations include a citation to the original paper.

Finally, we identify all cases in which a patent  $p$ , invented by inventors  $i$ , builds on a scientific idea that was first introduced in a paper co-authored by one of the inventors  $i$ . In these cases, the inventors should clearly be aware of the paper, having written it themselves. This results in 180,943 reusing patent–paper dyads. Even in this setting, where the paper is (co-)authored by the same inventor(s) and the idea is reused in a subsequent patent, the citation rate is only 36%, suggesting that in 64% of the cases the patent builds on the

scientific idea without citing the source paper.<sup>14</sup> To give one example, *adiponectin*—a protein hormone with antidiabetic and anti-inflammatory properties—was first discovered and published by Masashi Maeda, Iichiro Shimomura, and Yuji Matsuzawa. This scientific breakthrough transformed our understanding of fat tissue as an active endocrine organ and catalyzed a wave of scientific research on obesity, insulin resistance, and cardiovascular disease. It also spurred technological advances in diagnostics and therapeutics targeting metabolic disorders. The underlying scientific idea—captured by the noun phrase *adiponectin*—was subsequently reused in 23,688 scientific papers and 2,481 patents. Despite being co-authors of the original scientific discovery, the scientists did not cite their own paper in their first patent on adiponectin, filed a few years after the paper publication.

Taken together, these results show that even when there is strong indication that a patent draws directly on a specific scientific idea—and the inventor is arguably aware of the original paper—the patent often does not cite it. In line with prior survey evidence from inventors (Tijssen, 2002; Callaert et al., 2014), our findings suggest that patent-to-paper citations underestimate the extent to which patents rely on science. This pattern reinforces the view that patents do not necessarily cite all the scientific work on which they build, perhaps especially when the cited work is not legally required to delineate the scope of protection.

#### ***4.2 When do patents cite the scientific papers they build on?***

Building on our previous finding that patents often do not cite the scientific papers they rely on, we now examine the conditions under which patents are more likely to cite these source papers. Given that the primary objective of patent references is to delineate the legal scope of protection—rather than to acknowledge all scientific ideas upon which a patent builds—patents arguably prioritize citations to prior patents over scientific publications (Jaffe et al., 1993; Roach & Cohen, 2013). When patents do cite scientific papers, they are perhaps more common when the paper’s content is closely aligned with the patent, making it more relevant for defining legal boundaries, or when the paper is more recent, since newer work

---

<sup>14</sup> Remember that we are restricting the analysis to new ideas (noun phrases) that first appear in the scientific literature before being reused in patents. One might argue that inventors strategically avoid citing their own scientific work to reduce the risk of patent rejection. However, the grace period under U.S. patent law allows inventors to publicly disclose their invention—such as through a scientific publication or conference presentation—up to 12 months before filing a patent application without jeopardizing patentability. As a second strategy to avoid the risk that their own scientific publication becomes prior art that could block patentability, inventors may choose to file a patent application shortly before publishing the related scientific work. In such cases, the new idea (noun phrase) would typically appear first in a patent before entering the scientific literature. However, among all phrases first introduced in patents, only 1.7% subsequently first appear in a paper co-authored by at least one inventor of the pioneering patent, suggesting that only a very small minority are likely to represent patent–paper pairs.

better reflects the current state of the art while older scientific ideas may already be incorporated or cited by prior patents.

To test this idea, we use the same random sample of 100,000 scientific ideas (noun phrases) first introduced in scientific papers and later reused in at least one patent. For each idea, we identify all reusing patents, resulting in 7,139,906 patent–reused scientific idea dyads covering 2,468,508 unique patents. Using all these patent–reused scientific idea dyads, we then estimate the likelihood that the patent cites the original paper as a function of the time lag between the paper’s publication and the patent’s filing date. As shown in Figure 1, the predicted probability of citation declines steeply over time, falling to nearly zero for papers published more than 25 years earlier. Across the full population of U.S. patents, the scientific ideas reused within a patent are on average 46 years old at the time of filing, and in 75% of patents the average reused idea is at least 37 years old. This helps explain why many scientific papers are not cited, even when their ideas are incorporated into downstream patents.<sup>15</sup>

*‘Insert Figure 1’*

Next, we replicate the exact same analysis but predict the likelihood of citation as a function of the semantic instead of temporal distance between the patent and the paper on whose scientific idea(s) the patent builds. Specifically, we compute embedding vectors for each patent and paper based on their titles and abstracts, and calculate cosine similarity for all 7,139,906 patent–paper observations (Ghosh et al., 2024). As shown in Figure 2, the predicted probability of citation declines steeply as the semantic distance between the patent and the paper increases. This pattern further supports the idea that patents tend to cite scientific papers only when the underlying content is closely aligned and thus more likely to be viewed as legally relevant.<sup>16</sup>

*‘Insert Figure 2’*

---

<sup>15</sup> As previously discussed, besides directly citing the source paper, a patent may also cite it indirectly by referencing a paper or patent that cites the pioneering paper (e.g., Ahmadpoor & Jones, 2017). We find consistent results when using indirect citation as the outcome variable: both the predicted probability of citing a patent that cites the pioneering paper and the predicted probability of citing a paper that cites the pioneering paper decline sharply as the time lag between the patent’s filing and the pioneering paper’s publication increases. Results, not reported here to save space, are available from the authors upon request.

<sup>16</sup> We find consistent results when using indirect citation as the outcome: both the predicted probability of citing a patent that cites the pioneering paper and the predicted probability of citing a paper that cites the pioneering paper decrease sharply as the semantic distance between the focal patent and the source paper increases. Results, not reported here to save space, are available from the authors upon request.

## 5. Patents' Reliance on Scientific Ideas and Private Commercial Value

To demonstrate the value of our new approach and open dataset for strategy and innovation research, we revisit a long-standing research question: Are firm inventions that build on science associated with greater private value (Fleming & Sorenson, 2004; Krieger et al., 2024)? We follow the estimation strategy of Krieger et al. (2024) and rely on the KPSS measure of patent value, which captures the abnormal stock market returns experienced by a firm around the patent grant date (Kogan et al., 2017). To measure a patent's reliance on scientific ideas, we incorporate both the traditional metric (*Science References*) and our newly developed text-based metric (*Science Phrases*).

We follow Krieger et al. (2024) and sample 1,331,880 U.S. patents granted to 7,112 public U.S. firms between 1980 and 2010. A key feature and limitation of the KPSS data is that abnormal stock returns are observed at the firm-week level, rather than at the level of individual patents. Because firms often have multiple patents granted in the same week, these patents are all associated with a single stock market reaction, and each patent is assigned the average return for that firm-week. Following the preferred estimation strategy of Krieger et al. (2024), we also aggregate the independent variables to the firm-week level to match the unit of observation, computing their averages across all patents granted to the same firm within a given week. We then estimate weighted regressions at the firm-week level, using the number of patents granted as weights.

*'Insert Table 5'*

Table 5 reports results from specifications estimated without firm fixed effects (columns 1–4) and with firm fixed effects (columns 5–8). The inclusion of firm fixed effects accounts for persistent, time-invariant heterogeneity across firms, such as differences in industry, R&D strategy, or firms' overall science reliance. Across all models, we find that *Science Phrases* is positively and significantly associated with higher patent value. A one standard deviation increase in *Science Phrases* corresponds to a 22.3% increase in patent value without firm fixed effects (column 2) and an 10.9% increase with firm fixed effects (column 6). In line with prior work, *Science References* also show a positive association: a one standard deviation increase is linked to an 10.1% increase in patent value without firm fixed effects (column 1) and a 12.9% increase with firm fixed effects (column 5). Columns 4 and 8 re-estimate the correlation between *Science Phrases* and patent value for the subsample

of firm-weeks in which none of the firm's patents contain any science references.

Interestingly, the effect of *Science Phrases* remains positive, significant at  $p < 0.001$ , and large in magnitude in this restricted sample: a one standard deviation increase corresponds to a 23.1% increase in patent value without firm fixed effects (column 4) and a 8.1% increase with firm fixed effects (column 8). This suggests that *Science Phrases* captures remaining variation in a patent's reliance on science that is correlated with private patent value, even when no patent-to-paper citations are present.

We conduct several robustness checks. One potential concern is that, even though all regressions control for patent text length, longer patents may mechanically contain more *Science Phrases*, which could bias the measure. A related concern is that multiple distinct noun phrases within a patent may be synonyms and capture the same underlying scientific idea, artificially inflating *Science Phrases*. Both issues are mitigated by using *Science Phrases Share*, i.e. the proportion of phrases reused from the scientific literature relative to all unique noun phrases in the patent. As shown in Table A.5 in Appendix, our main results remain consistent when using *Science Phrases Share* instead of *Science Phrases*. Second, a key limitation is that abnormal stock returns are observed at the firm-week level, which provides a noisy proxy for the private value of individual patents when multiple patents are granted to the same firm in a given week. This concern is especially relevant because patent value is highly skewed and unlikely to be evenly distributed across patents within a firm-week. To address this issue, we followed the preferred specification of Krieger et al. (2024) and use the firm-week rather than the patent as our unit of analysis. Nevertheless, our results remain robust when using the patent as the unit of observation (Table A.6) and when restricting the sample to single-patent firm-weeks, where stock market reactions provide a cleaner signal of private patent value (Table A.7). Finally, private patent value based on stock market reactions is observable only for patents assigned to U.S. public firms and represents just one proxy for patent value. To broaden the analysis, Table A.8 replicates the results for the full population of U.S. patents, including those assigned to non-U.S. or privately held firms, using forward citations as an alternative measure of patent value (Harhoff et al., 1999).

Taken together, these results reaffirm that firm inventions building on science are more valuable (Fleming & Sorenson, 2004; Krieger et al., 2024). Our patent-level findings align with firm-level survey evidence showing that patent-to-paper citations underestimate the true impact of science on firm innovation performance (Roach & Cohen, 2013). By detecting

the reuse of scientific ideas even when uncited, our text-based approach complements citation-based indicators, broadens the lens on science–technology linkages, and uncovers firm performance effects that traditional measures overlook.

## **6. Discussion and Conclusion**

How do firms turn scientific knowledge into technological progress and competitive advantage? Strategy research has long emphasized that access to and absorption of scientific ideas is central to firm innovation. Yet despite this wide recognition by scholars, firms and policymakers, much of the science that underpins firm technologies leaves no trace in the conventional metric of patent-to-paper citations. This gap raises a central question: to what extent do firms rely on scientific ideas that remain invisible in traditional citation data, and does this hidden reliance matter for the value of firm innovation?

We address this question by developing a novel text-based measure that identifies when patents reuse specific scientific ideas from the academic literature. Unlike patent-to-paper citations, which primarily serve legal functions, this method captures the underlying scientific ideas on which patents are built, regardless of whether a citation is provided. In doing so, we open the black box of uncited science and trace its role in firm innovation. Our findings show that firms frequently build on scientific ideas without citing the original source paper—particularly when the science is older or less similar in content. Crucially, this hidden reliance is economically meaningful: patents that draw more heavily on scientific ideas are associated with higher private value, even after controlling for patent-to-paper citations and within the subset of patents without such citations. Put differently, science shapes firm innovation in ways that conventional citation indicators miss.

This study makes three contributions to strategy and innovation research. First, it develops and validates a scalable complementary measure of firms’ reliance on scientific ideas at the level of individual patents. Second, it shows that uncited science is systematically associated with private patent value, revealing a blind spot in citation-based approaches. Third, it introduces an open dataset that links the complete corpus of all U.S. patents with the corpus of all scientific publications until 2023. This resource enables researchers to trace idea-level connections between science and technology at scale—capturing idea-level flows from science to technology.

The open method and dataset we provide enables strategy and innovation scholars to revisit foundational questions and to address new ones about when, how, and why scientific ideas translate into firm innovation and competitive advantage. For instance: How do firms access and absorb external scientific knowledge in the process of innovation (Fleming & Sorenson, 2004)? How does reliance on science interact with firms' internal capabilities (Cohen & Levinthal, 1990)? And under what conditions does it translate into superior innovation performance (Krieger et al., 2024)? Importantly, our work also creates opportunities for future work to study at scale the two-way relationship in which science and technology co-evolve through iterative feedback loops (Rosenberg, 1976; Stokes, 1997; Barbosu & Teodoridis, 2025). While patent-to-paper citations are inherently unidirectional and thus only capture flows from science to technology, the text-based idea linkages in our dataset allow researchers for the first time to trace connections in both directions, from science to technology but also technology to science. Such bidirectional linkages cannot be observed using patent-to-paper citations, which are inherently unidirectional. In principle, one could study flows from technology to science via paper-to-patent citations, but such citations are very rare (Hsiao & Torvik, 2020). To give one example, *dendrimer* first appeared in a 1983 Dow Chemical patent describing a new class of perfectly branched, tree-like polymers, sparking an entire field of scientific research on nanoscale materials, polymer chemistry, and drug delivery, as evidenced by its reuse in more than 12,000 scientific papers.

Naturally, our study has limitations. Not all technologies are patented, and patenting propensities vary across industries and firms. Moreover, inventors and attorneys may strategically frame or draft patent texts—just as they may strategically cite prior art—which could affect how scientific ideas are reflected in our metric. Finally, firms may draw on science in ways not captured by either patent citations or our text-based measures. These limitations underscore that our approach should be viewed as a complement, rather than a substitute, for existing measures of science–technology linkages.

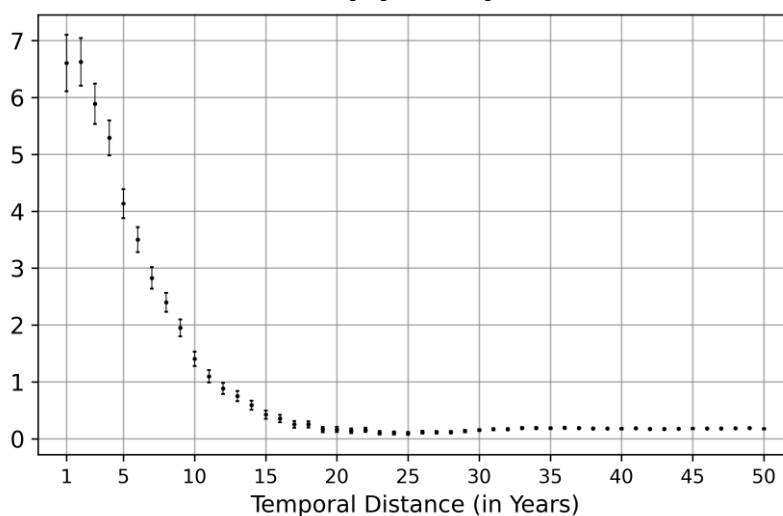
In sum, this paper shows that scientific ideas—cited or uncited—are a critical foundation of firm innovation and competitive advantage. By uncovering hidden reliance on science, our text-based approach expands the empirical toolkit available to strategy and innovation scholars and opens new avenues for understanding how science shapes firm performance.

## REFERENCES

- Ahmadpoor, M., & Jones, B. F. (2017). The dual frontier: Patented inventions and prior scientific advance. *Science*, 357(6351), 583-587.
- Arora, A., Belenzon, S., Cioaca, L., Sheer, L., Shin, H.M. & Shvadron, D. (2024). DISCERN 2.0: Duke Innovation & SCientific Enterprises Research Network [Dataset]. In Zenodo (CERN European Organization for Nuclear Research).  
<https://doi.org/10.5281/zenodo.3594642>
- Arts, S., Melluso, N., & Veugelers, R. (2025a). Beyond citations: Measuring novel scientific ideas and their impact in publication text. Forthcoming *Review of Economics and Statistics*. [https://doi.org/10.1162/rest\\_a\\_01561](https://doi.org/10.1162/rest_a_01561)
- Arts S, Fleming L, Veretennik L. (2025b). Harnessing academic science for corporate technology: The role of interpersonal networks and absorptive capacity.
- Balsmeier, B., Lück, S., & Fleming, L. (2025). Science knowledge localizes. *Research Policy*, 54(10), 105333.
- Barbosu, S., & Teodoridis, F. (2025). Catalysts of Discovery: How Technologies Are Shaping the Future of Science and Innovation. Cham, Switzerland: Palgrave Macmillan / Springer Nature. <https://doi.org/10.1007/978-3-031-98341-2>
- Bryan, K. A., Ozcan, Y., & Sampat, B. (2020). In-text patent citations: A user's guide. *Research Policy*, 49(4), 103946.
- Bryan, K. A., Ozcan, Y., & Sampat, B. (2024). *The Paper Trail of Knowledge, Revisited*. Working paper.
- Bush, V. (1945). *Science, the Endless Frontier: A Report to the President* (Washington: United States Government Printing Office, 1945).
- Callaert, J., Pellens, M., & Van Looy, B. (2014). Sources of inspiration? Making sense of scientific references in patents. *Scientometrics*, 98(3), 1617-1629.
- Cohen, W. M., & Levinthal, D. A. (1990). Absorptive capacity: A new perspective on learning and innovation. *Administrative science quarterly*, 35(1), 128-152.
- Cohen, W. M., Nelson, R. R., & Walsh, J. P. (2002). Links and impacts: the influence of public research on industrial R&D. *Management Science*, 48(1), 1-23.
- Fleming, L., & Sorenson, O. (2004). Science as a map in technological search. *Strategic Management Journal*, 25(8-9), 909-928.
- Ghosh, M., Erhardt, S., Rose, M. E., Buunk, E., & Harhoff, D. (2024). PaECTER: Patent-level representation learning using citation-informed transformers. *arXiv preprint arXiv:2402.19411*.
- Harhoff, D., Narin, F., Scherer, F. M., & Vopel, K. (1999). Citation frequency and the value of patented inventions. *Review of Economics and Statistics*, 81(3), 511-515.
- Hsiao, T. K., & Torvik, V. I. (2020). Technology footprints in scientific discovery: Citation contexts of paper-to-patent citations. *Proceedings of the Association for Information Science and Technology*, 57(1), e337.
- Iaria, A., Schwarz, C., & Waldinger, F. (2018). Frontier knowledge and scientific production: Evidence from the collapse of international science. *Quarterly Journal of Economics*, 133(2), 927-991.
- Jaffe, A. B., Trajtenberg, M., & Henderson, R. (1993). Geographic localization of knowledge spillovers as evidenced by patent citations. *Quarterly journal of Economics*, 108(3), 577-598.
- Jaffe, A. B., Trajtenberg, M., & Fogarty, M. S. (2000). Knowledge spillovers and patent citations: Evidence from a survey of inventors. *American Economic Review*, 90(2), 215-218.

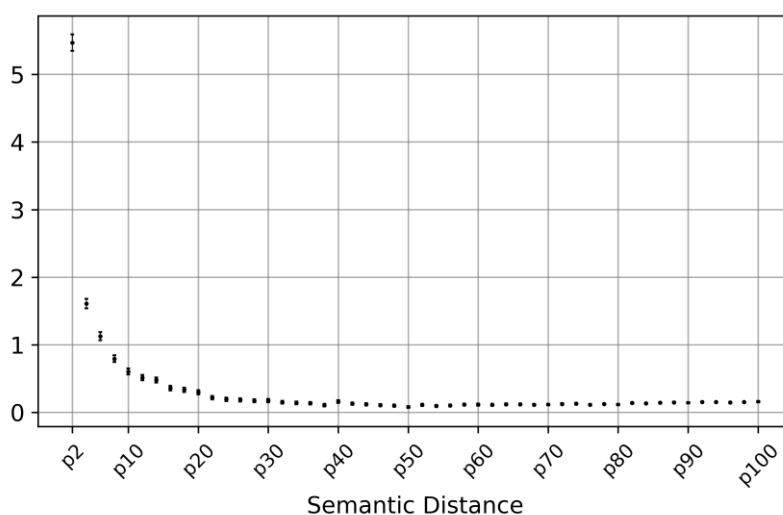
- Kogan, L., Papanikolaou, D., Seru, A., & Stoffman, N. (2017). Technological innovation, resource allocation, and growth. *Quarterly journal of Economics*, 132(2), 665-712.
- Krieger, J. L., Schnitzer, M., & Watzinger, M. (2024). Standing on the shoulders of science. *Strategic Management Journal*, 45(9), 1670-1695.
- Lampe, R. (2012). Strategic citation. *Review of Economics and Statistics*, 94(1), 320-333.
- Mansfield, E. (1991). Academic research and industrial innovation. *Research Policy*, 20(1), 1-12.
- Marx, M., & Fuegi, A. (2020). Reliance on science: Worldwide front-page patent citations to scientific articles. *Strategic Management Journal*, 41(9), 1572-1594.
- Mokyr, J. (2011). The gifts of Athena: Historical origins of the knowledge economy. In *The gifts of Athena*. Princeton University Press.
- Narin, F., Hamilton, K. S., & Olivastro, D. (1997). The increasing linkage between US technology and public science. *Research Policy*, 26(3), 317-330.
- Roach, M., & Cohen, W. M. (2013). Lens or prism? Patent citations as a measure of knowledge flows from public research. *Management Science*, 59(2), 504-525.
- Rosenberg, N. (1976). Technological change in the machine tool industry, 1840–1910. In *Perspectives on Technology* (pp. 9–31). Cambridge University Press.
- Scharfmann, E., Marx, M., & Fleming, L. (2024). Scientists invent better patents, and inventors publish better science. Working paper.
- Stephan, P. (2012). *How economics shapes science*. Harvard University Press.
- Stokes, D. E. (1997). *Pasteur's quadrant: Basic science and technological innovation*. Brookings Institution Press.
- Teplitskiy, M., Duede, E., Menietti, M., & Lakhani, K. R. (2022). How status of research papers affects the way they are read and cited. *Research policy*, 51(4), 104484.
- Tijssen, R. J. (2002). Science dependence of technologies: evidence from inventions and their inventors. *Research policy*, 31(4), 509-526.
- Verluisse, C., Cristelli, G., Higham, K., & de Rassenfosse, G. (2025). Beyond the front page: In-text citations to patents as traces of inventor knowledge. Forthcoming *Strategic Management Journal*. <https://doi.org/10.1002/smj.70027>

**Figure 1: Predicted probability that patents cite the scientific papers they build on, by time lag between paper and patent**



*Notes:* This figure is based on 7,139,906 patent–scientific idea (noun phrase) dyads corresponding to 2,468,508 unique patents each reusing at least one of 100,000 randomly sampled scientific ideas (noun phrases) first introduced in the scientific literature. It plots the predicted probability (in %) that a patent cites the pioneering paper that introduced the idea it reuses, as a function of the time lag (in years) between the paper’s publication and the patent’s filing. Predictions and 95% confidence intervals are based on a linear probability model with year-lag specific indicators. The model includes fixed effects for the scientific subfield (OpenAlex 4-digits) of the pioneering paper, the technology subclass (CPC 4-digits) of the reusing patent, and the publication year of the paper. Additional controls for the pioneering paper include abstract availability, text length (number of unique phrases in title and abstract), and the number of papers cited. Controls for the reusing patent include text length (number of unique phrases in the full text), and the number of cited scientific papers.

**Figure 2: Predicted probability that patents cite the scientific papers they build on, by content similarity between paper and patent**



*Notes:* This figure is based on 7,139,906 patent–scientific idea (noun phrase) dyads corresponding to 2,468,508 unique patents each reusing at least one of 100,000 randomly sampled scientific ideas (noun phrases) first introduced in the scientific literature. It plots the predicted probability (in %) that a patent cites the pioneering paper that introduced the idea it reuses, as a function of the semantic distance between the reusing patent and the pioneering paper, modeled as a series of fifty 2-percentile bins. Bin p2 represents patents most similar in content to the pioneering paper, while bin p100 represents patents most distant in content. Predictions and 95% confidence intervals are based on a linear probability model. The model includes fixed effects for the subfield (OpenAlex 4-digits) of the pioneering paper, the technology subclass (CPC 4-digits) of the reusing patent, and the publication year of the paper. Additional controls for the pioneering paper include abstract availability, text length (number of unique phrases in title and abstract), and the number of papers and journals cited. Controls for the reusing patent include text length (number of unique phrases in the full text), and the number of cited patents and scientific papers.

**Table 1: Patent-Level Summary Statistics**

	Full Sample (n = 7,699,586)						
	Mean	Std	Min	p25	p50	p75	Max
Phrases	254.663	307.310	1.000	118.000	179.000	281.000	81,913.000
Science Phrases Bin	0.963	0.188	0.000	1.000	1.000	1.000	1.000
Science Phrases	32.232	87.023	0.000	4.000	9.000	24.000	10,435.000
Science Phrases Share	0.081	0.078	0.000	0.029	0.053	0.105	0.643
Science References Bin	0.221	0.415	0.000	0.000	0.000	0.000	1.000
Science References	2.913	17.469	0.000	0.000	0.000	0.000	1,502.000
Science Phrases Avg Age	45.823	13.189	0.000	36.923	45.800	54.677	120.000
Science References Avg Age	10.627	8.577	0.000	5.000	9.000	14.000	202.000

Notes: US patents granted between 1976 and June 2023. *Science Phrases Avg Age* and *Science References Avg Age* are calculated only for patents with at least one reused scientific noun phrase and at least one scientific reference, respectively.

**Table 2: Summary Statistics Across Patents with Scientific References, University Assignees, and Scientist Inventors**

	Mean	Std	Min	Max	Mean	Std	Min	Max	Z	p-value
	Patents with Scientific References (n=1,785,379)				Patents with no Scientific References (n= 5,914,207)					
Science Phrases Bin	0.995	0.067	0.000	1.000	0.953	0.211	0.000	1.000	85.217	0.000
Science Phrases	91.549	160.586	0.000	10,435.000	14.325	26.299	0.000	5,330.000	1,258.985	0.000
Science Phrases Share	0.148	0.101	0.000	0.632	0.061	0.056	0.000	0.643	1,131.971	0.000
Science References Bin	1.000	0.000	1.000	1.000	0.000	0.000	0.000	0.000	2,028.338	0.000
Science References	12.814	34.604	1.000	1,502.000	0.000	0.000	0.000	0.000	2,028.338	0.000
	University Patents (n=321,086)				Non-University Patents (n=7,378,500)					
Science Phrases Bin	0.987	0.115	0.000	1.000	0.962	0.191	0.000	1.000	23.394	0.000
Science Phrases	102.842	164.626	0.000	4,535.000	29.159	80.603	0.000	10,435.000	433.945	0.000
Science Phrases Share	0.161	0.109	0.000	0.627	0.078	0.074	0.000	0.643	451.529	0.000
Science References Bin	0.627	0.484	0.000	1.000	0.215	0.411	0.000	1.000	395.757	0.000
Science References	16.151	39.774	0.000	964.000	2.398	15.606	0.000	1,502.000	446.967	0.000
	Patents with Scientist Inventors (n=3,945,226)				Patents without Scientist Inventors (n=3,754,360)					
Science Phrases Bin	0.981	0.136	0.000	1.000	0.944	0.229	0.000	1.000	88.739	0.000
Science Phrases	48.610	112.739	0.000	10,435.000	15.022	39.962	0.000	5,465.000	891.668	0.000
Science Phrases Share	0.103	0.088	0.000	0.643	0.058	0.056	0.000	0.624	786.205	0.000
Science References Bin	0.328	0.469	0.000	1.000	0.131	0.338	0.000	1.000	472.384	0.000
Science References	4.926	23.052	0.000	1,502.000	0.917	7.925	0.000	1,380.000	496.787	0.000

Notes: Sample includes 7,699,586 U.S. patents granted between 1976 and 2023. “Patents with Scientific References” are those citing at least one scientific publication, the comparison group is all other patents. “University Patents” are those with at least one university assignee, the comparison group is all other patents. “Patents with Scientist Inventors” are those with at least one inventor who has authored a scientific publication (Scharfmann et al., 2024), the comparison group is all others. Z-values are based on Mann–Whitney tests assuming unequal variances

**Table 3: (Indirect) Citation, Geographic, and Social Distance between Inventors and Authors**

	Mean (%)	Std (%)	Mean (%)	Std (%)	Z	p-value
Panel A	reusing patent–paper dyads (n=7,139,906)		matched-control patent–paper dyads (n=7,139,906)			
patent-to-paper citation	0.340	5.825	0.005	0.712	10.9743	0.000
indirect patent-to-paper citation (via cited patent)	0.301	5.482	0.015	1.235	9.3651	0.000
indirect patent-to-paper citation (via cited paper)	1.252	11.119	0.075	2.729	38.5341	0.000
inventor–author geographic distance	7.123	1.736	7.342	1.625	-178.306	0.000
inventor–author overlap (social distance = 0)	0.097	3.116	0.002	0.410	3.125	0.000
direct inventor–author tie (social distance = 1)	0.101	3.169	0.011	1.046	2.931	0.000
indirect tie via common collaborator (social distance = 2)	0.553	7.412	0.222	4.701	10.830	0.000
no tie within two degrees (social distance ≥ 3)	99.250	8.628	99.766	4.833	-16.886	0.000
Panel B	reusing patent–paper dyads time lag ≤5years (n=105,219)		matched-control patent–paper dyads time lag ≤5years (n=105,219)			
patent-to-paper citation	6.367	24.416	0.087	2.940	24.950	0.000
indirect patent-to-paper citation (via cited patent)	1.624	12.641	0.038	1.949	6.302	0.000
indirect patent-to-paper citation (via cited paper)	7.705	26.667	0.180	4.234	29.896	0.000
inventor–author geographic distance	7.103	1.720	7.347	1.616	-24.841	0.000
inventor–author overlap (social distance = 0)	2.743	16.336	0.038	1.949	10.749	0.000
direct inventor–author tie (social distance = 1)	1.947	13.818	0.105	3.232	7.321	0.000
indirect tie via common collaborator (social distance = 2)	7.676	26.622	2.976	16.992	18.675	0.000
no tie within two degrees (social distance ≥ 3)	87.633	32.921	96.882	17.381	-36.745	0.000

*Notes:* The unit of observation is a patent–paper dyad constructed from a random sample of 100,000 scientific ideas (noun phrases) first introduced in scientific papers. The sample includes 7,139,906 reusing patent–paper dyads linking the source paper that first introduced a scientific idea to a patent that reuses that idea, and an equal number of matched-control dyads linking the same source paper to a patent that does not reuse the idea but is filed and granted in the same years and assigned to the same technology subclass. *Patent-to-paper citation* equals one if the patent cites the source paper. *Indirect citation (via cited patent)* equals one if the patent cites another patent that cites the source paper; *indirect citation (via cited paper)* equals one if the patent cites a scientific paper that cites the source paper. *Inventor–author geographic distance* is the shortest distance, in miles, between any inventor location on the patent and any author affiliation on the source paper, calculated from PatentsView and OpenAlex coordinates (missing for 21.9% of papers and 0.4% of patents). *Social distance* is the shortest interpersonal path between any inventor and author based on collaboration networks observed up to the filing year, connecting all OpenAlex authors and U.S. patent inventors through prior co-authorships and co-inventorships. Inventor–author linkages rely on the disambiguation and crosswalk by Scharfmann et al. (2024). *Inventor–author overlap (social distance = 0)* equals one if an inventor is also an author of the source paper; *direct inventor–author tie (social distance = 1)* if they have co-authored or co-invented before filing; *indirect tie via common collaborator (social distance = 2)* if connected through a shared collaborator; *no tie within two degrees (social distance ≥ 3)* otherwise. All distances are measured as of the filing year. Binary variables are expressed as percentages. *Inventor–author geographic distance* is measured in miles and log-transformed after adding one. Z-values are based on Mann–Whitney tests assuming unequal variances.

**Table 4: Ordered Logit Model Linking the Share of a Firm’s R&D Lab Projects That Used Public Research to Patents’ Reliance on Science**

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Science References Avg	0.887 (0.160) [0.000]		0.788 (0.164) [0.000]		0.626 (0.232) [0.007]		0.523 (0.237) [0.027]	
Science Phrases Avg		0.508 (0.140) [0.000]	0.354 (0.144) [0.014]	0.371 (0.187) [0.047]		0.514 (0.207) [0.013]	0.411 (0.209) [0.049]	0.414 (0.252) [0.101]
Observations	646	646	646	349	646	646	646	349
Pseudo R2	0.0396	0.0282	0.0440	0.0094	0.0879	0.0870	0.0913	0.0954
Log-Likelihood	-691.3	-699.5	-688.1	-333.9	-656.5	-657.2	-654.1	-305
Industry Fixed Effects	No	No	No	No	Yes	Yes	Yes	Yes
Sample	Full	Full	Full	No Science Ref.	Full	Full	Full	No Science Ref.

*Notes:* The unit of observation is an R&D lab surveyed by the 1994 Carnegie Mellon Survey on Industrial R&D with at least one patent filed in the three years prior to the survey year. Ordered logit regressions link the survey-based measure of the share of a firm’s R&D unit’s projects that used public research to patent-based indicators of reliance on scientific ideas. The outcome variable is the ordered response to the following question of the Carnegie Mellon Survey on Industrial R&D: “During the last three years, what percentage of your R&D unit’s projects made use of research findings produced by universities or government research labs?” with five categories (Below 10%, 10–40%, 41–60%, 61–90%, Over 90%). *Science References Avg* is the average logged number of scientific references per patent filed by the R&D lab in the three years preceding the survey, and *Science Phrases Avg* is the average logged number of scientific noun phrases per patent filed by the R&D lab in the three years preceding the survey. All regressions control for the average text length of the patents (average logged number of noun phrases) and the logged number of patents filed by the R&D lab in the three years preceding the survey. Columns (1)–(4) report specifications without industry fixed effects; columns (5)–(8) include industry fixed effects. Columns (4) and (8) restrict the sample to R&D labs for which none of the patents filed by the lab in the previous three years contained any scientific references. Robust standard errors are reported in parentheses, and p-values in brackets.

**Table 5: Patents' Reliance on Scientific Ideas and Private Value**

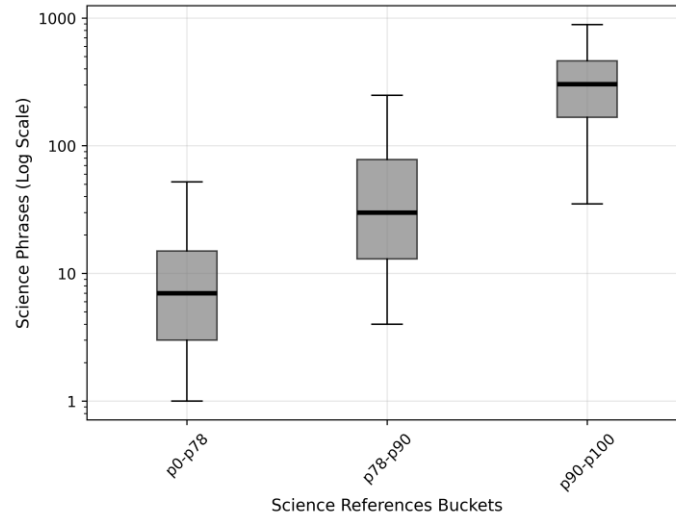
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Science References	2.050 (1.025) [0.046]		1.295 (1.032) [0.209]		2.617 (0.755) [0.001]		2.428 (0.740) [0.001]	
Science Phrases		3.179 (1.091) [0.004]	2.797 (1.092) [0.010]	3.722 (0.693) [0.000]		1.566 (0.509) [0.002]	1.089 (0.466) [0.019]	1.302 (0.348) [0.000]
Observations	394,546	394,546	394,546	241,446	393,196	393,196	393,196	240,037
R-squared	0.114	0.114	0.115	0.082	0.386	0.385	0.386	0.430
Technology Controls	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Week Fixed Effects	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Firm Fixed Effects	No	No	No	No	Yes	Yes	Yes	Yes
Sample	Full	Full	Full	No Science Ref.	Full	Full	Full	No Science Ref.
Marginal Effects (in %)								
Science References	10.127		6.397		12.898		11.968	
Science Phrases		22.248	19.574	23.125		10.932	7.598	8.054

*Notes:* OLS regressions link patent value (KPSS) to reliance on scientific ideas as measured by the number of *Science References* and *Science Phrases*. The unit of observation is the firm-week. The dependent variable is the mean KPSS of all patents granted to a firm in a given week. Independent variables are log-transformed at the patent level and averaged within the firm-week. Regressions are weighted by the number of patents in the firm-week. All models include issue-week and control for the full vector of 4-digit CPC technology classes by including, for each CPC class, the average share of the firm-week's patents assigned to that class. Controls also include the text length of the patent (average number of noun phrases). The sample includes 1,331,880 U.S. patents granted to 7,112 public U.S. firms between 1980 and 2010. Models 4 and 8 focus only on firm-weeks in which none of the firm's patents contain science references. Noun phrases are extracted from the full patent text. Marginal effects represent the percentage change in average KPSS per firm-week for a one-standard deviation increase in the independent variable. Robust standard errors (clustered at the firm level) are reported in parentheses. P-values of the estimated coefficients are reported in brackets.

## ONLINE APPENDIX

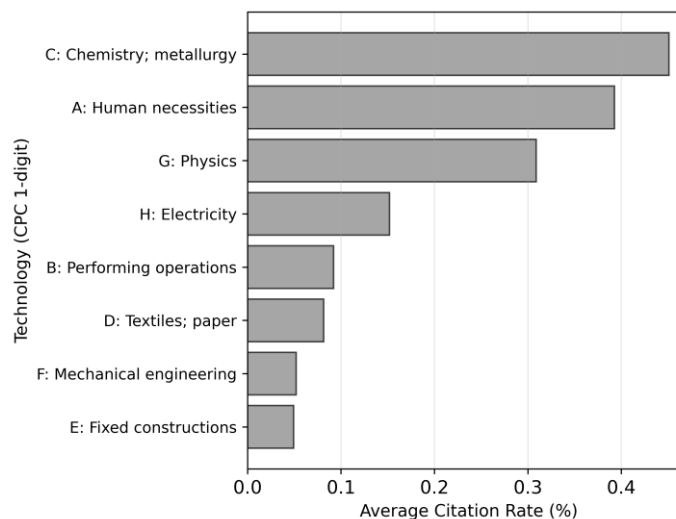
### Appendix A

**Figure A.1 Distribution of *Science Phrases* across Percentiles of *Science References***



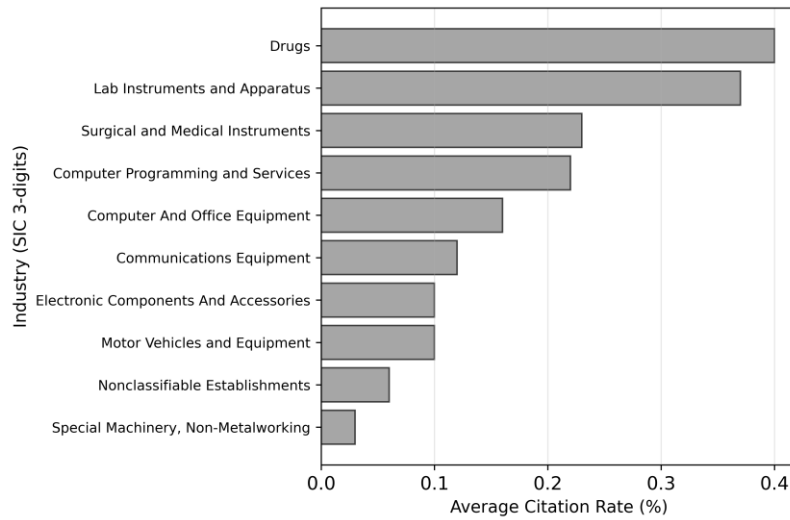
*Notes:* This figure shows box plots of the distribution of *Science Phrases* across percentile buckets of *Science References* for all 7,699,586 U.S. patents granted between 1976 and 2023. The x-axis groups patents into three buckets of *Science References*: p0–p78, p78–p90, and p90–p100. Since 78% of patents cite no scientific paper, the first bucket captures this large share with zero patent-to-paper citations. The y-axis uses a logarithmic scale for visualization, but the tick labels are reported in original counts for ease of interpretation. Boxes represent the interquartile range (25th–75th percentile), the thick horizontal line marks the median, and whiskers extend from the 5th to 95th percentiles.

**Figure A.2 Average Citation Rate by Technology**



*Notes:* This figure reports the average citation rate across technologies (CPC 1-digit) for 7,139,906 patent–phrase dyads. These dyads are constructed from a random sample of 100,000 scientific ideas (noun phrases) first introduced in the scientific literature and identifying all patents reusing these phrases in their text. The citation rate (in %) corresponds to the share of patent–phrase dyads in which the patent cites the pioneering paper where the reused phrase first appeared.

**Figure A.3 Average Citation Rate by Industry**



*Notes:* This figure reports the average citation rate across industries (SIC 3-digit) for 2,217,751 patent–phrase dyads. These dyads are constructed from a random sample of 100,000 scientific ideas (noun phrases) first introduced in the scientific literature and identifying all patents filed by U.S. public firms that reuse these phrases. The sample covers 851,385 granted patents filed between 1980 and 2021 by 6,088 unique firms, as linked through DISCERN2.0. The citation rate (in %) corresponds to the share of patent–phrase dyads in which the patent cites the pioneering paper where the reused phrase first appeared. Industries with zero citation rates or fewer than 50 firms are excluded. Industry names longer than 30 characters are abbreviated using GPT-4o-mini.

**Table A.1 Most Reused Scientific Phrases in Patents, by Decade**

1951-1960	1961-1970	1971-1980	1981-1990
copolymer (365,333)	polyethylene_glycol (260,386)	liquid_crystal_display (432,075)	wide_area_network (420,415)
polypropylene (326,747)	silicone (229,688)	optical_storage (241,886)	field-programmable_gate_array (142,338)
polyester (324,224)	peptide (229,436)	ethernet (237,520)	electroporation (79,755)
monomer (295,132)	polycarbonate (222,169)	gsm (175,930)	metropolitan_area_network (62,824)
polyurethane (244,726)	polyamide (207,337)	liquid_crystal (146,299)	apoptosis (61,471)
isopropanol (166,842)	chemical_vapor_deposition (195,345)	magnetic_disk_storage (144,070)	protocol_stack (49,610)
silica_gel (155,672)	tetrahydrofuran (190,691)	amino_acid_sequence (137,754)	transcription_factor (46,880)
styrene (151,477)	multiplexer (169,995)	plasmid (128,924)	solid-state_drive (45,399)
cathode_ray_tube (151,142)	epoxy_resin (168,313)	nucleotide_sequence (127,637)	transgene (45,116)
elastomer (142,313)	polypeptide (166,255)	dna_sequence (123,808)	interleukin (44,898)
1991-2000	2001-2010	2011-2020	
hypertext_transfer_protocol (74,748)	microrna (16,593)	smart_contract (4,400)	
extensible_markup_language (64,998)	bevacizumab (13,701)	nivolumab (4,130)	
atomic_layer_deposition (61,256)	imatinib (11,805)	pembrolizumab (3,649)	
carbon_nanotube (45,643)	cetuximab (11,666)	cryptocurrency (2,938)	
organic_light-emit_diode (37,992)	erlotinib (10,408)	generative_adversarial_network (2,883)	
paclitaxel (32,555)	gefitinib (10,386)	beidou_navigation_satellite (2,789)	
graphene (29,191)	bortezomib (9,176)	cas9 (2,710)	
nand_flash_memory (26,987)	everolimus (8,138)	vemurafenib (2,548)	
phage_display (26,460)	mimo_transmission (7,885)	tensorflow (2,500)	
qr_code (25,913)	abiotic_stress_resistance (7,855)	kubernetete (2,228)	

*Notes:* Top 10 noun phrases originating in scientific publications and subsequently reused in patents, grouped by the decade of first adoption in patents. The decade is defined by the filing year of the earliest patent containing the phrase. Values in parentheses indicate the total number of distinct patents in which the phrase appears through 2023.

**Table A.2 Most Reused Scientific Phrases by Selected Industries and Firms**

Industry	Firm	N. Patents	Science Phrases (avg)	Top Reused Science Phrases
Chemicals and Allied Products	Genentech	1,715	607	polypeptide (90.9%); amino_acid_sequence (85.8%); peptide (81.6%); plasmid (79.9%); dna_sequence (79.1%)
	Incyte	1,206	426	bind_affinity (85.6%); radionuclide (83.7%); hydrophobicity (82.5%); magnesium_stearate (81.9%); polyvinylpyrrolidone (81.3%)
	Abbvie	1,685	424	polyethylene_glycol (58.6%); liposome (46.0%); peptide (42.6%); magnesium_stearate (40.1%); cycloalkyl (38.1%)
	Amgen	2,021	415	peptide (73.2%); polypeptide (68.7%); polyethylene_glycol (67.7%); amino_acid_sequence (67.6%); monoclonal_antibody (59.8%)
	Ionis Pharm	1,240	372	oligonucleotide (97.3%); phosphorothioate (93.2%); nucleoside (89.5%); mma (87.3%); cytosine (87.1%)
Measuring, Analyzing, and Controlling Instruments	Bio Rad Lab	1,751	147	peptide (41.5%); oligonucleotide (24.9%); polyethylene_glycol (23.9%); dichloromethane (23.8%); solvate (23.6%)
	Hologic	1,028	106	nucleic_acid_sequence (32.4%); oligonucleotide (30.7%); amplification_reaction (29.2%); nucleic_acid_amplification (27.5%); polynucleotide (25.8%)
	Abbott Lab	9,374	104	polyurethane (21.6%); polyester (20.0%); copolymer (20.0%); silicone (19.9%); polyethylene_glycol (19.0%)
	Thermo Fisher	3,853	104	peptide (24.9%); oligonucleotide (23.1%); mass_spectrometry (21.7%); mass_spectrometer (21.5%); ion_source (17.7%)
	Waters	1,201	68	mass_spectrometer (68.2%); mass_spectrometry (61.6%); liquid_chromatography (59.2%); ion_source (50.9%); mass_analyser (42.3%)
Petroleum Refining and Related Industries	Exxon Mobil	11,746	81	copolymer (32.0%); monomer (27.2%); polyolefin (20.9%); molecular_sieve (20.8%); polypropylene (20.4%)
	Chevron	6,212	65	copolymer (21.4%); monomer (20.0%); molecular_sieve (19.2%); cyclohexane (16.3%); n-hexane (15.9%)
	Amoco	2,160	52	monomer (19.3%); copolymer (19.2%); catalyst_bed (14.7%); molecular_sieve (13.3%); polyester (13.3%)
	Mobil	4,888	42	silica-alumina (27.6%); extrudate (22.5%); x-ray_diffraction (20.4%); ion_exchange (20.4%); molecular_sieve (17.7%)
	Texaco	2,770	33	synthesis_gas (12.5%); polyol (11.6%); copolymer (11.2%); alkylene_oxide (9.2%); pore_volume (8.9%)
Oil and Gas Extraction	Conocophillips	4,924	40	copolymer (24.2%); monomer (23.0%); cyclohexane (18.7%); homopolymer (14.5%); pore_volume (13.9%)
	Shell Oil	1,085	34	copolymer (21.8%); cyclohexanone (15.3%); tetrahydrofuran (13.6%); corrosion_inhibitor (13.3%); aqueous_dispersions (12.6%)
	Halliburton	11,673	25	proppant (17.2%); gravel_pack (13.1%); copolymer (12.2%); cement_slurry (10.2%); fiber_optic (10.1%)
	Schlumberger	10,283	23	proppant (7.6%); gamma_ray (7.4%); seismic_data (7.1%); elastomer (5.9%); gravel_pack (5.7%)
	Baker Intl	7,652	20	tungsten_carbide (26.0%); polycrystalline_diamond (21.5%); downhole_motor (18.9%); hydrocarbon_production (18.0%); drillstre (16.5%)

Notes: For each selected industry (SIC 2-digit code), the table lists the top five U.S. public firms ranked by the average number of distinct scientific noun phrases reused per patent. For each firm, the five most frequently reused phrases are shown, with percentages indicating the share of the firm's patents in which the phrase appears. For example, Genentech averages 607 scientific phrases per patent, and 90.9% of its patents contain the phrase polypeptide. Firm–patent linkages are based on DISCERN2.0. The sample includes U.S. granted patents from 1980–2021 and is restricted to firms with at least 1,000 patents.

**Table A.3: Summary Statistics Restricted to Patents without Scientific References**

	Mean	Std	Min	Max	Mean	Std	Min	Max	Z	p-value
	University Patents (n=119,888)				Non-University Patents (n=5,794,319)					
Science Phrases Bin	0.968	0.176	0.000	1.000	0.953	0.211	0.000	1.000	8.751	0.000
Science Phrases	19.209	29.715	0.000	1,004.000	14.224	26.214	0.000	5,330.000	94.303	0.000
Science Phrases Share	0.082	0.067	0.000	0.627	0.061	0.055	0.000	0.643	123.860	0.000
	Patents with Scientist Inventors (n=2,652,135)				Patents without Scientist Inventors (n= 3,262,072)					
Science Phrases Bin	0.973	0.162	0.000	1.000	0.938	0.242	0.000	1.000	74.595	0.000
Science Phrases	19.382	32.869	0.000	2,822.000	10.215	18.382	0.000	5,330.000	585.799	0.000
Science Phrases Share	0.074	0.063	0.000	0.643	0.051	0.046	0.000	0.624	482.371	0.000

*Notes:* Sample restricted to 5,914,207 U.S. patents granted between 1976 and 2023 with no citations to papers. “University Patents” are those with at least one university assignee, the comparison group is all others. “Patents with Scientist Inventors” are those with at least one inventor who has authored a scientific publication (Scharfmann et al., 2024), the comparison group is all others. Z-values are from Mann-Whitney tests.

**Table A.4: OLS Model Linking R&D Labs’ Reliance on Open Science to Patents’ Reliance on Scientific Ideas**

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Science References Avg	0.474		0.417		0.297		0.257	
	(0.068)		(0.071)		(0.097)		(0.098)	
	[0.000]		[0.000]		[0.002]		[0.009]	
Science Phrases Avg		0.268	0.182	0.189		0.198	0.152	0.140
		(0.067)	(0.069)	(0.087)		(0.091)	(0.090)	(0.108)
		[0.000]	[0.008]	[0.030]		[0.030]	[0.093]	[0.198]
Observations	646	646	646	349	646	646	646	349
R-squared	0.091	0.067	0.101	0.019	0.209	0.205	0.213	0.152
Industry Fixed Effects	No	No	No	No	Yes	Yes	Yes	Yes
Sample	Full	Full	Full	No Science Ref.	Full	Full	Full	No Science Ref.

*Notes:* The unit of observation is an R&D lab surveyed by the 1994 Carnegie Mellon Survey on Industrial R&D with at least one patent filed in the three years prior to the survey year. OLS regressions link the “open science” factor from Roach and Cohen (2013) to patent-based indicators of reliance on scientific ideas. The outcome variable is the continuous factor score constructed from R&D managers’ responses rating, on a four-point Likert scale, the importance of publications and reports, public conferences and meetings, and informal information exchange as channels through which public research contributes to their R&D. *Science References Avg* is the average logged number of scientific references per patent filed by the R&D lab in the three years preceding the survey, and *Science Phrases Avg* is the average logged number of scientific noun phrases per patent filed by the lab in the same period. All regressions control for the average text length of the patents (average logged number of noun phrases) and the logged number of patents filed by the lab during 1991–1993. Columns (1)–(4) report specifications without industry fixed effects; columns (5)–(8) include industry fixed effects. Columns (4) and (8) restrict the sample to R&D labs for which none of the patents filed by the lab in the previous three years contained any scientific references. Robust standard errors are reported in parentheses, and p-values in brackets.

**Table A.5: Patents' Reliance on Scientific Ideas and Private Value: *Science Phrases Share***

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Science References	2.050 (1.025) [0.046]		1.272 (1.025) [0.215]		2.617 (0.755) [0.001]		2.271 (0.693) [0.001]	
Science Phrases Share		28.095 (14.673) [0.056]	22.745 (15.109) [0.132]	35.987 (10.013) [0.000]		23.537 (7.304) [0.001]	16.712 (6.163) [0.007]	13.925 (4.678) [0.003]
Observations	394,546	394,546	394,546	241,446	393,196	393,196	393,196	240,037
R-squared	0.114	0.114	0.114	0.081	0.386	0.385	0.386	0.430
Technology Controls	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Week Fixed Effects	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Firm Fixed Effects	No	No	No	No	Yes	Yes	Yes	Yes
Sample	Full	Full	Full	No Science Ref.	Full	Full	Full	No Science Ref.
Marginal Effects (in %)								
Science References	10.127		6.282		12.898		11.194	
Science Phrases		12.921	10.460	13.345		10.797	7.666	5.140

*Notes:* OLS regressions link patent value (KPSS) to reliance on scientific ideas as measured by the number of *Science References* and *Science Phrases Share*. The unit of observation is the firm-week. The dependent variable is the mean KPSS of all patents granted to a firm in a given week. Independent variables are log-transformed at the patent level and averaged within the firm-week. Regressions are weighted by the number of patents in the firm-week. All models include issue-week and control for the full vector of 4-digit CPC technology classes by including, for each CPC class, the average share of the firm-week's patents assigned to that class. Controls also include the text length of the patent (average number of noun phrases). The sample includes 1,331,880 U.S. patents granted to 7,112 public U.S. firms between 1980 and 2010. Models 4 and 8 include only firm-weeks where all patents have no science references. Noun phrases are extracted from the full patent text. Marginal effects represent the percentage change in average KPSS per firm-week for a one-standard deviation increase in the independent variable. Robust standard errors (clustered at the firm level) are reported in parentheses. P-values of the estimated coefficients are reported in brackets.

**Table A.6: Patents' Reliance on Scientific Ideas and Private Value: Patent-Level Analysis**

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Science References	1.717 (0.496) [0.001]		1.290 (0.487) [0.008]		1.176 (0.313) [0.000]		1.067 (0.315) [0.001]	
Science Phrases		2.831 (0.481) [0.000]	2.578 (0.470) [0.000]	2.436 (0.415) [0.000]		0.897 (0.165) [0.000]	0.716 (0.158) [0.000]	0.646 (0.132) [0.000]
Observations	1,331,870	1,331,870	1,331,870	1,033,650	1,330,532	1,330,532	1,330,532	1,032,258
R-squared	0.086	0.087	0.087	0.079	0.357	0.357	0.357	0.343
Technology Fixed Effects	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Week Fixed Effects	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Firm Fixed Effects	No	No	No	No	Yes	Yes	Yes	Yes
Sample	Full	Full	Full	No Science Ref.	Full	Full	Full	No Science Ref.
Marginal Effects (in %)								
Science References	10.599		7.964		7.256		6.583	
Science Phrases		27.997	25.496	23.627		8.863	7.078	6.259

*Notes:* OLS regressions link patent value (KPSS) to reliance on scientific ideas as measured by the number of *Science References* and *Science Phrases*. The unit of observation is a single patent. The dependent variable is the KPSS. Independent variables are log-transformed. All models include issue-week and 4-digit CPC technology-class fixed effects, and control for the text length of the patent (number of noun phrases). The sample includes 1,331,880 U.S. patents granted to 7,112 public U.S. firms between 1980 and 2010. Models 4 and 8 include only patents with no science references. Noun phrases are extracted from the full patent text. Marginal effects represent the percentage change in KPSS for a one-standard deviation increase in the independent variable. Robust standard errors (clustered at the firm level) are reported in parentheses. P-values of the estimated coefficients are reported in brackets.

**Table A.7: Patents' Reliance on Scientific Ideas and Private Value: Patent-Level Analysis for Single-Patent Firm-Weeks**

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Science References	-0.397 (0.706) [.574]		-1.132 (0.713) [.113]		1.863 (0.488) [.000]		1.757 (0.490) [.000]	
Science Phrases		4.581 (0.728) [.000]	4.855 (0.734) [.000]	4.196 (0.644) [.000]		1.143 (0.249) [.000]	0.834 (0.242) [.001]	0.782 (0.248) [.002]
Observations	218,954	218,954	218,954	163,215	217,594	217,594	217,594	161,790
R-squared	0.073	0.075	0.075	0.080	0.569	0.569	0.569	0.579
Technology Fixed Effects	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Week Fixed Effects	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Firm Fixed Effects	No	No	No	No	Yes	Yes	Yes	Yes
Sample	Single-patent weeks	Single-patent weeks	Single-patent weeks	Single-patent weeks with no science ref.	Single-patent weeks	Single-patent weeks	Single-patent weeks	Single-patent weeks with no science ref.
Marginal Effects (in %)								
Science References	-2.067		-5.897		9.668		9.120	
Science Phrases		32.201	34.124	25.049		7.997	5.837	4.639

*Notes:* OLS regressions link patent value (KPSS) to reliance on scientific ideas as measured by the number of *Science References* and *Science Phrases*. The unit of observation is the firm-week with only one patent. The dependent variable is KPSS of the single patent granted to a firm in a given week. Independent variables are log-transformed at the patent level and averaged within the firm-week. Regressions are weighted by the number of patents in the firm-week. All models include issue-week and 4-digit CPC technology-class fixed effects, and control for the text length of the patent (average number of noun phrases). The sample includes patents granted to 7,093 public U.S. firms between 1980 and 2010. Noun phrases are extracted from the full patent text. Marginal effects represent the percentage change in the KPSS per patent for a one-standard deviation increase in the independent variable. Standard errors are clustered at the firm level.

**Table A.8: Patents' Reliance on Scientific Ideas and Forward Citations**

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Science References	0.120 (0.002) [.000]		0.106 (0.002) [.000]		0.106 (0.004) [.000]		0.099 (0.004) [.000]	
Science Phrases		0.107 (0.002) [.000]	0.076 (0.002) [.000]	0.101 (0.002) [.000]		0.071 (0.004) [.000]	0.047 (0.004) [.000]	0.068 (0.004) [.000]
Observations	4,138,107	4,138,107	4,138,107	3,313,476	1,504,485	1,504,485	1,504,485	1,160,256
Pseudo R2	0.220	0.218	0.220	0.197	0.275	0.273	0.275	0.255
Year FE	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Tech FE	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Firm FE	No	No	No	No	Yes	Yes	Yes	Yes
Sample	Full	Full	Full	No Science Ref.	Full	Full	Full	No Science Ref.
Marginal Effects								
Science References	9.363		8.264		8.141		7.594	
Science Phrases		13.811	9.760	10.904		8.762	5.820	7.099

*Notes:* Poisson quasi-maximum likelihood regressions link forward citations to reliance on scientific ideas, measured by the number of *Science References* and *Science Phrases*. The unit of observation is the patent. The dependent variable is the number of forward citations within 10 years of grant. Independent variables are log-transformed after adding one. All models include grant year and technology class (4-digit CPC) fixed effects, and control for patent text length (number of noun phrases). Models 1–4 include all U.S. patents granted between 1980 and 2012; Models 5–8 restrict to patents from 7,303 unique U.S. public firms; Models 4 and 8 further restrict to patents with no science references. Noun phrases are extracted from the full patent text. Marginal effects represent the percentage change in forward citations for a one-standard-deviation increase in the explanatory variable.

## Appendix B - New Phrase Scoring

To identify scientific ideas, we build on Arts et al. (2025a), who processed the full text of all scientific publications and conference proceedings in OpenAlex from 1666 to 2023, yielding 27,079,343 distinct noun phrases. We then apply a large-scale classification procedure using the GPT-4o-mini model, which assigns each phrase a confidence score indicating the likelihood that it reflects a genuine scientific idea. We submitted batch requests via the OpenAI API. Each batch contained 100 noun phrases presented in a structured input table, with each row listing the phrase, the title of the first paper in which it appeared, the journal name, and the year of publication. The model returned an idea score on a 1–3 scale:

- 1 = Very unlikely to represent a new idea
- 2 = Likely to represent a new idea
- 3 = Very likely to represent a new idea

The exact prompt used for classification was as follows:

```
# Input Data
You are given a data table where each record contains the following fields:
- phrase: A noun phrase extracted from the title or abstract of a scientific paper.
- title: The title of the paper.
- journal: The name of the journal in which the paper was published.
- year: The publication year of the paper.

# Background Context
- The noun phrases have been extracted from all papers indexed in OpenAlex.
- Each noun phrase is linked to the first corresponding paper in which it appeared (by publication date) within the entire corpus of indexed papers.

# Task
Evaluate whether the given noun phrase represents a new scientific idea introduced by the specified paper in the given year. A "new scientific idea" should be interpreted broadly and may include:
- A novel concept, theory, method, tool, technology, or discovery (e.g., a newly identified molecule or chemical structure).

Your judgment should be based on all available information provided in the dataset but particularly on all other available information at your disposal.

# Response Format
Provide your evaluation using the following structure:
- Idea Score: Assign a score from 1 to 3 indicating how likely the noun phrase represents a new scientific idea introduced by the given paper and year.
-- 1: Very unlikely
-- 3: Very likely

Ensure your assessment is well-reasoned and considers the provided metadata and all other available information at your disposal in your judgment.
```

Across 27,079,601 unique new phrases, the model assigned:

- 2,403,306 phrases (8.9%) a score of 1
- 13,771,184 phrases (50.8%) a score of 2
- 10,905,111 phrases (40.3%) a score of 3

We construct two alternative versions of the *New Phrase* measure by applying cutoffs to the idea scores:

- *New Phrase (Score 2 and 3)*: excludes phrases with a score of 1
- *New Phrase (Score 3)*: retains only phrases scored as 3, excluding both 1 and 2

Tables B.1 and B.2 below provide summary statistics at the phrase and paper level respectively.

**Table B.1: Summary Statistics Phrase Level**

	#	Mean	Std	Min	Reuse in subsequent papers					Max	Skew
					p25	p50	p75	p95			
New Phrase (All Scores)	27,079,601	17.228	434.835	1	1	2	5	34	416,734	315.421	
New Phrase (Score 2 and 3)	24,676,295	15.503	396.043	1	1	2	5	31	416,734	346.605	
New Phrase (Score 3)	10,905,111	18.357	445.352	1	1	2	5	36	344,357	289.012	

*Notes:* *New Phrase (All Scores)* includes all new phrases with no exclusions. *New Phrase (Score 2 and 3)* includes only phrases with score 2 and 3, thus excluding phrases with score 1. *New Phrase (Score 3)* includes only phrases Score 3, thus excluding phrases with score 1 and 2. Summary statistics of the reuse of new phrases in subsequent papers, i.e. how many later papers reuse the phrase in their text. # is the number of units. p25, p50, p75, p95 are respectively the 25th, 50th, 75th, 95th percentile. The skewness (skew) of the distributions is the measure of the asymmetry of the probability distribution of a real-valued random variable about its mean. A positive skewness indicates that the distribution has a long right tail, while a negative skewness indicates a long left tail.

**Table B.2: Summary Statistics Paper Level**

	Mean	Std. Dev	Min	p25	p50	p75	p95	Max	Skew
New Phrase (All Scores) (Binary)	0.245	0.430	0	0	0	0	1	1	1.184
New Phrase (All Scores)	0.360	0.779	0	0	0	0	2	223	4.128
New Phrase Reuse (All Scores)	6.555	268.222	0	0	0	0	12	416,735	522.944
New Phrase (Score 2 and 3) (Binary)	0.229	0.420	0	0	0	0	1	1	1.292
New Phrase (Score 2 and 3)	0.328	0.735	0	0	0	0	2	216	4.262
New Phrase Reuse (Score 2 and 3)	5.081	233.776	0	0	0	0	9	416,734	609.017
New Phrase (Score 3) (Binary)	0.111	0.314	0	0	0	0	1	1	2.476
New Phrase (Score 3)	0.145	0.471	0	0	0	0	1	43	4.735
New Phrase Reuse (Score 3)	2.659	173.650	0	0	0	0	3	344,357	768.870

*Notes:* *New Phrase (All Scores)* includes all new phrases with no exclusions. *New Phrase (Score 2 and 3)* includes only phrases with score 2 and 3, thus excluding phrases with score 1. *New Phrase (Score 3)* includes only phrases Score 3, thus excluding phrases with score 1 and 2. n= 75,295,921 papers published between 1901 and 2023. p25, p50, p75, p95 are respectively the 25th, 50th, 75th, 95th percentile. The skewness (skew) of the distributions is the measure of the asymmetry of the probability distribution of a real-valued random variable about its mean. A positive skewness indicates that the distribution has a long right tail, while a negative skewness indicates a long left tail.

To assess whether the GPT-4o-mini model produces a meaningful confidence score—that is, whether higher-scoring phrases more accurately capture genuine scientific ideas—we replicate the validation exercise from Arts et al. (2025a)<sup>17</sup> on distinguishing Nobel Prize papers from matched controls (results displayed in Table B.3 below).

**Table B.3: Nobel Prize Classification**

Variable	Precision	Recall	AUC	Marginal Effect
Panel A: Ex-ante				
New Phrase (All Scores)	69.87	65.63	0.7635	20.93
New Phrase (Score 2 and 3)	71.35	66.15	0.7738	21.90
New Phrase (Score 3)	73.98	66.15	0.7923	24.60
Panel B: Ex-post				
New Phrase Reuse (All Scores)	74.42	67.19	0.7938	25.66
New Phrase Reuse (Score 2 and 3)	74.85	67.19	0.7980	25.84
New Phrase Reuse (Score 3)	76.32	67.71	0.8063	30.99

*Notes:* Reproduction of Table 3 Column 2 (Panel A of current table) and Column 11 (Panel B of current table) from Arts et al. (2025a). *New Phrase (All Scores)* includes all new phrases with no exclusions. *New Phrase (Score 2 and 3)* includes only phrases with score 2 and 3, thus excluding phrases with score 1. *New Phrase (Score 3)* includes only phrases Score 3, thus excluding phrases with score 1 and 2. Logit regressions with a binary outcome for Nobel Prize paper, robust standard errors in parentheses. Sample includes 1,168 papers (584 Nobel Prize papers matched with 584 control papers from the same year, journal, and subfield). All measures are log-transformed (after adding 1 for zero values). Models control for publication year and subfield fixed effects, abstract availability, text length (unique words and phrases in title and abstract), and number of unique papers and journals cited. AUC represents the area under the ROC curve. Marginal effects show the percentage increase in the likelihood of being a Nobel Prize paper associated with a one-standard-deviation increase in the metric.

<sup>17</sup> The replication data and code of Arts et al. (2025) are retrieved from <https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/OUNQTO>

The results show that applying stricter cutoffs significantly improves predictive validity. For Nobel Prize papers, precision<sup>18</sup> rises from 69.9% for *New Phrase (All Scores)*, to 71.4% for *New Phrase (Score 2 and 3)*, and to 74.0% for *New Phrase (Score 3)*, with the AUC increasing from 0.76 to 0.77 and 0.79 respectively. As expected, excluding phrases with a score of 1 primarily filters out false positives and improves precision without substantially affecting false negatives or recall. Weighting noun phrases by their reuse in subsequent papers—capturing the impact and influence of new ideas on future scientific discourse—strengthens these effects further as illustrated in Panel B of Table B.3. In addition, the marginal effect—calculated as the percentage increase in the probability that a paper is a Nobel Prize paper associated with a one-standard-deviation increase in the respective metric—grows steadily across cutoffs, from 20.9% for *New Phrase (All Scores)* to 21.9% for *New Phrase (Score 2 and 3)*, and to 24.6% for *New Phrase (Score 3)*. When weighting phrases by their reuse in subsequent papers, marginal effects become even larger, reaching 31.0% for *New Phrase Reuse (Score 3)*.

---

<sup>18</sup> Precision is the proportion of true positives among all predicted positives (i.e., the share of papers classified as Nobel Prize papers that are actually Nobel Prize papers). Recall is the proportion of true positives among all actual positives (i.e., the share of Nobel Prize papers correctly identified). AUC (area under the ROC curve) summarizes overall predictive accuracy by capturing the trade-off between true positives and false positives across thresholds.