

The Private Value of Open-Source Innovation*

Logan P. Emery[†] Chan Lim[‡] Shiwei Ye[†]

December 2024

Abstract

We investigate open-source innovation by public firms and the private value it generates for these firms. Unlike patents, which grant inventors exclusive rights to their inventions, open-source innovations can be used by anyone. Nevertheless, using an extensive dataset of public-firm activity on GitHub, we find that firms with open-source projects represent 68% of the U.S. stock market across 86% of industries. We estimate the private value of all projects in our sample to be nearly \$25 billion, with the average project generating \$842,000. We find that projects with fully permissive licenses are generally less valuable and firms facing higher competition tend to generate less private value from their projects. We also find that complementarity with commercial products is not a primary driver of private value. Finally, open-source value significantly predicts firm growth in terms of sales, profits, employment, and patenting. These results contribute to our understanding of the private value generated by innovation in the absence of excludability.

Keywords: Open Source, Innovation, Firm growth, Valuation, GitHub
JEL Classification: G14, G30, L21, O36

*We thank Fabrizio Corre, Mikael Paaso, and seminar participants at Erasmus University Rotterdam and Monash University for helpful comments. We are responsible for any remaining errors. Emails: emery@rsm.nl, clim26@buffalo.edu, ye@rsm.nl.

[†]Rotterdam School of Management, Erasmus University Rotterdam.

[‡]School of Management, University at Buffalo.

1 Introduction

Innovation is costly to produce but can be used by multiple parties with little or no additional cost. Consequently, inventors bear the full cost of innovation but receive only a portion of the benefits. In order to incentivize investment in innovation, systems (e.g., patents) have been created to grant inventors exclusive rights to monetize their innovation for a given period of time. This excludability is seen as crucial for deriving private value from innovation (Arrow, 1962; Crouzet et al., 2022). However, the past decade has seen a rise in so-called open-source innovation. When an innovation is “open-sourced,” it is made publicly available to all parties at little or no cost. Perhaps surprisingly, many firms choose to make their innovation open source. A recent survey finds that 90% of Fortune 100 companies use GitHub, the largest platform for developing open-source innovation.¹ However, it remains unclear what private value these profit-maximizing entities derive from making their costly innovation freely available.

In this paper, we study open-source innovation by publicly traded firms. Our analysis proceeds in three parts. First, we document the extent to which publicly traded firms produce open-source innovation and characterize the type of firm that chooses to develop their innovation via open source. Second, we estimate the private value of open-source innovation and investigate the innovation, firm, and product market characteristics that most strongly correlate with private value. Finally, we examine the relation between open-source innovation and future firm growth.

There are many ways in which open-source innovation may create private value for firms. We group these ways into three broad mechanisms.² First, making innovation free-to-use

¹ See <https://octoverse.github.com/2022/>.

² These mechanisms are drawn from a broad literature that considers the possible incentives for freely revealing innovation. These papers include Allen (1983), Lerner and Tirole (2002), Harhoff et al. (2003), Dahlander and Gann (2010), Henkel et al. (2014), Parker et al. (2017), Alexy et al. (2018), Nagle (2018), Teece (2018) and Lin and Maruping (2022). For reviews of the open-source literature, see von Hippel and von Krogh (2003), Goldfarb and Tucker (2019), and Dahlander et al. (2021).

can maximize the adoption of that innovation. Increased adoption can create value by improving the innovation through community development, giving the firm more control over subsequent development of the technology, creating a product ecosystem that deters customers from switching to competitors, and increasing demand for complementary commercial products. These phenomena constitute network effects whereby the value of the innovation increases exponentially with the number of users. Second, open-source innovation may provide value through labor considerations. Valuable employees may want to share their work to increase their reputation within the open-source community, and firms may identify talented workers who contribute to open-source projects (and have lower integration costs due to already being familiar with the firm’s infrastructure). Finally, making innovation open source may enhance the firm’s reputation. Firms that contribute to open-source projects may be seen as more community oriented and open-source projects may be seen as more transparent and certified for quality.

We study open-source activity using public-firm activity on GitHub. While not all open-source innovation takes place on GitHub, it is the largest platform for developing open-source innovation, specifically computer software, and has become synonymous with the idea of open source. We compile a comprehensive dataset of public-firm activity on GitHub from 2015 through 2023. While only 18% of public firms produce open-source innovation on GitHub (“open-source firms”), those firms represent 68% of the total stock market capitalization and 80% of the total research and development (R&D) expenditure by public firms in 2023. In comparison, open-source firms represented only 20% of the total stock market capitalization in 2015. Moreover, while 32% of open-source firms are from the “Computer Software” industry, 86% of industries have at least one such firm,³ demonstrating the growing scope of open-source innovation. We find that firms participating in the open source innovation are larger, more valuable, more innovative, and face less competition on average. However, in a regression setting, most of these differences are absorbed by firm fixed effects, suggesting

³ This analysis is based on the Fama-French 49 industries classification system.

that firm fixed effects can account for much of the potential selection bias in which firms choose to make innovation open source.

We next employ a modified version of the method developed by [Kogan et al. \(2017\)](#) to measure the private value of repositories as estimated by investors. The methodology is based on observing firm-specific stock returns over the three days following the release of a repository. The resulting estimates reflect the value captured by the firm (i.e., private value), excluding any value generated for other firms (i.e., public value), and represent the sum of the value of the innovation and the value of being open source. Using this methodology, we find that investors estimate the average repository in our sample to generate \$842,849 (in 2023 dollars), with average values increasing significantly over the sample period. The total private value created by repositories in our sample is nearly \$25 billion. The most valuable GitHub portfolios are owned by Amazon.com, Inc and Microsoft Corp, both valued at nearly \$8 billion, and repositories using Python as the main programming language produce the most value in our sample.

Attributing stock returns around the repository announcement to the repository assumes that investors respond to the release of repositories on GitHub. Otherwise, the estimates of repository value only reflect noise in the market. To assess this possibility, we regress repository value on a measure of future repository popularity (i.e., stars). We find that more-valuable repositories end up being significantly more popular in the future, which suggests that stock price reactions to repository announcements contain value-relevant information. To provide further validation, we perform a placebo test where we assign random announcement days, within the true announcement year, to each repository and estimate a placebo value. We then perform the same regression of (placebo) repository value on future repository popularity. In none of the 500 iterations of this test does the placebo relation reach the economic significance of the true relation.

Investigating the determinants of open-source value, we find that repositories with copy-left licenses, which place restrictions on commercial use of the repository, are more valuable

on average than repositories with fully permissive licenses. This result highlights the value of excludability even in an open-source setting. We also find evidence against complementarity between open-source projects and commercial products being a first-order driver of open-source value. Instead, standalone open-source projects tend to be more valuable and represent a larger fraction of firms' open-source-portfolio value on average. Larger repositories (e.g., more lines of code) are also not necessarily more valuable and repositories with more subsequent issues opened (e.g., bugs) are perceived as less valuable when initially released.

Since a firm's competitors are most likely to benefit from the open-source nature of GitHub repositories, we also investigate how product market characteristics correlate with the private value of open-source innovation. We find that firms facing less competition tend to have repositories that produce more private value. This is potentially a function of these firms capturing a larger portion of the total value created by the repository, which is the sum of the private and public values. Alternatively, these firms may be more willing to share valuable innovation due to fewer concerns of negative competitive consequences. In either case, we find that competition is a significant consideration for firms producing open-source innovation. We also find that, controlling for the level of competition, firms that are more likely to benefit from spillover effects in the product market produce more-valuable repositories, which is consistent with the importance of network effects for open-source value.

Finally, we investigate the relation between open-source innovation and future firm growth. We find that firms generating more open-source value have a larger growth in sales, profits, number of employees, and both the number and value of patents granted over the following three years. Thus, open-source innovation produces significant value for the innovator despite it being available for use by competitors.

This paper makes several contributions to the literature on innovation. First, our paper contributes to the broad literature on measuring the economic value of innovation. Existing studies have typically explored the value of innovation within traditional intellectual property

protection systems, such as patents or trademarks, which grant exclusive rights to use and monetize innovative outputs (Pakes, 1985; Austin, 1993; Hall et al., 2005; Kogan et al., 2017; Chen et al., 2019; Desai et al., 2023; Ahmadi et al., 2024). We explore how innovative outputs contribute to a company’s value even when freely disclosed to a broad audience through open-source licenses. Specifically, in the case of open-source software, as noted by Lerner and Tirole (2005a), the contribution of a company’s intellectual assets to its value creation and future growth is often indirect, making it challenging to measure quantitatively. We address this challenge by leveraging financial markets to measure the value of intellectual property without excludability through the value of repositories in open-source platforms.

Most directly, our paper contributes to the literature on open-source innovation. We construct an extensive dataset of open-source activity by public firms on GitHub, which allows us to document open-source activity at a granular level. Furthermore, it allows us to develop a new stock-market-based measure of the value of open-source innovation. Previous research has estimated the value of particular open-source software, such as Apache and nginx, using the cost of replicating similar services with proprietary software (Greenstein and Nagle, 2014; Murciano-Goroff et al., 2021). Other research estimates the aggregate economic value of open-source software using a cost-to-produce approach, based on the length of software code and labor costs (Robbins et al., 2021; Blind et al., 2021). For example, Hoffmann et al. (2024) use this strategy to estimate the cost of replicating the most-used open-source software either once (\$4.15 billion) or individually by all firms that use it (\$8.8 trillion).⁴ In contrast, we measure the private value of open-source activity by public firms using stock-market reactions, allowing us to quantify the dollar value of individual repositories at the firm level and explore heterogeneity therein.

We also contribute to the literature on innovation and firm growth by showing that the value of open-source innovation provides significant insights into firm growth beyond what is captured by the measure of patent value proposed by Kogan et al. (2017). Previous

⁴ Outside of the realm of open source, Gómez-Cram and Lawrence (2024) investigate the value of software by studying the long-run stock returns of software companies.

research has shown that companies can enhance their software development capabilities, firm productivity, and access to venture capital by being active in open-source platforms (Nagle, 2018, 2019; Conti et al., 2021). Our study adds to this literature by testing the impact of open-source innovation on a company’s long-term growth in sales, profits, employment, and innovative output.

2 Institutional Background

In this section, we discuss the process of developing open-source innovation on the GitHub platform and provide institutional details necessary for understanding the data used in our analysis.

2.1 Initiating open source projects

To deploy their projects on GitHub, firms need to create organization accounts. While some firms create only one organization account, others create multiple organization accounts based on organizational divisions, purposes, or related products. Repositories (projects) can then be created within these accounts, and administrators decide whether the projects will be publicly visible or only visible privately to certain organization or project members with the necessary permissions. The creation and management of public repositories come with almost no costs, while support and some features for managing private repositories require GitHub Team or GitHub Enterprise subscriptions. Notably, even though there had been additional costs for adding private repositories before the GitHub pricing model was changed from repository-based to user-based in 2015, GitHub has provided free hosting for public repositories since its inception.

One crucial decision to make when creating a repository is choosing a license. Without a license, projects cannot be considered open source, even if the source code is publicly visible.⁵

⁵ <https://choosealicense.com/no-permission/>

The choice of license can have different implications for commercial use. There are two primary categories of open-source licenses based on their permissiveness: permissive licenses and copyleft licenses. Permissive licenses impose minimal restrictions on how the source code can be used. In contrast, copyleft licenses require that (part of) derivative projects using the licensed code must also be open source. Therefore, firms that intend to find a balance between sharing their work with the community and protecting their proprietary interests may find copyleft licenses more attractive, as their competitors may be hesitant to open source their proprietary developments built upon copyleft-licensed projects. It is worth noting that some firms also opt for customized licenses with clauses that effectively limit commercial use. These custom licenses may appear open source but include restrictions that make the projects more “source available” rather than truly open source.⁶ [Internet Appendix A](#) compares different types of licenses based on permissions and conditions.

2.2 Project development and community interaction

GitHub operates on the Git system, a collaborative and distributed platform for software development. In this context, several key processes and community interactions play a crucial role in fostering innovation and progress. This section provides a brief overview of these processes.

The development process begins with the creation of a repository, where developers work on the code locally on their own computers. Changes are saved using the “commit” command, which records updates to the local repository along with brief summaries describing the modifications. Each commit serves as a checkpoint, documenting what was done and why. When developers are ready, they “push” or upload these commits to the remote repository, making the updates accessible to other contributors and users. This workflow enables efficient collaboration and version control, driving the iterative improvement of software projects.

⁶ <https://opensource.org/osd/>

Users interested in staying updated on a repository’s progress can “star” a repository, essentially bookmarking it for future reference. Those who have questions or suggestions can also “open issues.” Both the development team and fellow community members actively participate in addressing these issues.

Furthermore, users can engage in the development process by “forking” the repository, which allows them to create a personal copy and work on the codebase independently. If the changes made in this personal fork are deemed valuable and applicable to the original project, users can initiate “pull requests.” These pull requests serve as formal requests to integrate the changes back into the original repository. The changes proposed in pull requests undergo review and, if approved, are merged into the main codebase, thereby contributing to the open-source project’s ongoing development.

3 Open-Source Activity

3.1 Data

To construct our dataset of GitHub activities by U.S. public firms, we begin by linking GitHub organization accounts with firms. Following the methodology of [Conti et al. \(2021\)](#), we first collect websites of organization accounts via the GHTorrent project and the GitHub API. We then compare these domains with the web URLs of U.S. public firms and their subsidiaries from Compustat or Orbis. To ensure the accuracy of our matches, we screen out accounts whose domains are indicative of hosting or social media services, such as “github.com” and “facebook.com.” We then conduct a rigorous manual search to complement our domain-based matching. Specifically, we query the firm names together with the term “open source” via Google to locate official web pages that list their open source projects, and search the firm names on GitHub to identify associated organization accounts. Following this, we compile a comprehensive list of public repositories tied to the identified organization accounts through the GHArchive database, which records and archives times-

tamped public activity of GitHub repositories. In total, we match 1,281 firms with 3,314 organization accounts and 168,085 public repositories up to the year 2023.

Upon establishing a link between U.S. public firms and their respective GitHub organization accounts and public repositories, we utilize the GHArchive to gather additional information on the public footprints of these repositories. Most importantly, we determine the dates when the repositories were made public by identifying timestamps associated with the earliest activity, specifically those labeled as “PublicEvent.” Pinpointing the exact dates is crucial for our valuation process, which ultimately depends on the stock market reaction. We also create a firm-month panel that includes measures of aggregated activities observable to the public, such as the cumulative counts of repositories and the number of opened issues, pushes, and pull requests. Our panel spans the years 2015 to 2023, representing a relatively comprehensive picture of organizational engagement within the open-source community.

Additionally, we employ the GitHub API to collect static characteristics of 140,824 repositories extant as of February 2024. This includes an array of attributes from descriptive repository metadata, such as creation dates, licenses, and programming languages, to quantitative measures of community engagement, including the number of stars, watchers, and forks.

Finally, we use large language models (LLMs) to classify or evaluate repositories based on topics, complementarity, and novelty. We use OpenAI’s API to interact with the GPT-4o model, providing information including the repository name, description, main programming language, self-reported topics, website, and the name of the repository owner. We then prompt the model to conduct evaluation tasks. For topics, we use the model to assign a relatedness score (0 to 1) to 17 pre-defined topic categories, constructed from the GitRanking taxonomy (Sas et al., 2023). We define the complementarity score (0 to 1) as the extent to which a repository complements the firm’s commercial products (instead of being a standalone product). We also use the model to evaluate the novelty of a repository (0 to 1), which measures how novel or groundbreaking it is compared to existing solutions, focusing

on whether it introduces new ideas, techniques, or approaches. We take various steps to ensure consistency across repositories, including clear definitions of evaluation tasks, scoring reference systems, and consistency checks by conducting multiple rounds of scoring on a small subsample to verify stable outputs for each repository. [Internet Appendix B](#) provides details of our approach, model parameters, and prompts used.

3.2 Summary statistics

Before moving to the estimation of open-source value, we first provide an overview of open-source activity. We begin by documenting trends in open-source engagement during our sample period, which are plotted in [Figure 1](#). The dashed yellow line plots the cumulative number of repositories created by firms that are public as of that month, which totals 122,107 repositories by the end of our sample period.⁷ The blue line represents the percentage of public firms that have created at least one repository on GitHub as of that month (henceforth “open-source firms”). This percentage increases steadily from 4.8% in January 2015 to 18.1% in December 2023. The red line plots the percentage of total market capitalization represented by open-source firms, which is 67.5% at the end of our sample period. Finally, the green line plots the percentage of total R&D expenditure represented by open-source firms, which is 80.2% at the end of our sample period. Thus, despite open-source firms being only one-fifth of all public firms, they represent two-thirds of the stock market and over four-fifths of investment in innovation. We therefore conclude that firms engaged in open-source activities are an important part of the US economy.

Next, we examine the distribution of open-source firms across industries. [Figure 2](#) presents pie charts of this distribution at the firm and repository levels. We use the Fama-French five industry classification scheme ([Fama and French, 1997](#)) and further separate out

⁷ Note that the cumulative counts of repositories in our firm-month panel are slightly smaller than the original matched sample for two reasons. First, some publicly visible repositories that never appear in major open-source event records (such as issues, pushes, and pull requests) are excluded from the panel. Second, we exclude delisted firms starting from the month of delisting.

the “Computer Software” and “Finance” industries as defined by the Fama-French 49 industry classification scheme. One may assume that software firms represent the majority of open-source firms in our sample. However, we find that only 32.2% of open-source firms in our sample come from the “Computer Software” industry. This said, more than two-thirds of repositories in our sample are owned by firms in this industry, confirming the intuition that software firms are the most active in open-source innovation. Nonetheless, other industries are also reasonably well represented in our sample, particularly the “Business Equipment, Telephone and Television Transmission” and “Consumer Durables, NonDurables, Wholesale, Retail, and Some Services” industries. The breadth of open-source engagement across industries likely reflects the growing importance of software across all parts of the economy.

We provide further summary statistics of open-source activities in Table 1. Panel A reports the distribution of repositories for all firms, as well as by industry, as of December 2023. As can be inferred from Figure 2, open-source activity is most common among “Computer Software” firms, with 62.6% of firms in this industry having open-source activity. The “Healthcare, Medical Equipment, and Drugs” industry is the least active, with only 6.4% of firms having open-source activity. This result may be a reflection of the importance of excludability for innovation in this industry.⁸

Panel B compares firm and product market characteristics for firms with and without open-source activity. The panel reports the mean and median values of each group of firms over our sample period. Firm characteristics include the number of employees, market-to-book ratio, return-on-assets, investment, sales growth, tangibility, and research and development (R&D) expenditure scaled by total assets, all of which are calculated using data from Compustat. We also calculate market capitalization and annual returns using data from the Center for Research in Security Prices (CRSP) and obtain data on patent portfolios from Kogan et al. (2017).

We also examine the product market characteristics of open-source firms because tech-

⁸ E.g., see <https://www.reuters.com/legal/litigation/us-senators-ask-regulators-clear-drug-patent-thickets-2022-06-08/>.

nological spillovers to competitors represent the largest potential negative externality faced by open-source firms. The product market characteristics we consider include market power, scope, product market centrality, product market similarity, and product market fluidity. Market power measures a firm’s dominance within its product market, which we proxy for using a structural estimate of markups from [Pellegrino \(2024\)](#). Scope measures the number of industries in which the firm operates, based on their product descriptions in SEC filings ([Hoberg and Phillips, 2024](#)). Product market centrality is calculated as the eigenvector centrality of a firm in the product-market network, which is constructed using similarity scores from [Hoberg and Phillips \(2016\)](#). This quantity reflects the extent of competition faced by the firm, but it also measures the extent to which the firm benefits from spillover effects in the network (i.e., network effects), which can be crucial for the success of open-source projects. Product market similarity measures how similar a firm’s products are to its peers’ ([Hoberg and Phillips, 2016](#)). Finally, product market fluidity measures how intensively a firm’s product market is changing ([Hoberg et al., 2014](#)). A description of each variable and its data source is provided in [Table A1](#).

We find that open-source firms are considerably larger than non-open-source firms on average, based on market capitalization, employees, and number of patents. These firms also tend to have higher valuations, based on market-to-book ratio, which could reflect investors’ assessment of growth opportunities resulting from the firms’ innovation. Intuitively, open-source firms tend to have less tangible assets and larger R&D expenditures. Finally, open-source firms appear to face less competition: they charge higher markups, have lower product market centrality, are less similar to their product market rivals, and operate in less fluid product markets.

While these summary statistics paint a preliminary picture of open-source firms, they do not account for the concentration of open-source activity in certain industries, particularly the “Computer Software” industry. The observed differences between open-source and non-open-source firms could, therefore, be a function of industry differences rather than firm-

specific characteristics. We investigate this possibility in the next section.

3.3 Determinants of open-source activity

To provide a more rigorous characterization of open-source firms, we next consider the determinants of open-source activity in a regression setting. This approach controls for time-varying industry and time-invariant firm fixed effects. It also allows us to test the relative strength of the correlation between variables and open-source activity to identify key determinants of open-source activity. These tests are intended to be descriptive rather than causal, helping researchers understand potential selection bias in open-source activity and identify relevant omitted variables in a given research context.

The results from this analysis are reported in Table 2. The data pertains to firm-month observations of all public firms. Each regression includes the full set of firm and product market characteristics discussed in the previous section, as well as industry-time fixed effects. Columns (2) and (4) report regressions that also include firm fixed effects. Standard errors are double clustered on industry and time, and all independent variables are standardized to facilitate interpretation.

The regressions reported in Columns (1) and (2) pertain to all public firms. The dependent variable is an indicator that equals one if the firm is an open-source firm.⁹ The results in Column (1) therefore reflect the average difference between open-source and non-open-source firms within each industry. We find that firms with higher valuations (Mkt Cap, Market-to-Book) and more innovation (N Patents, R&D Exp) are more likely to be open-source firms. However, these firms also appear to be less profitable than their industry peers (Return-on-Assets) and have lower annual returns. One possible interpretation is that open-source firms may follow a traditional strategy in technology industries of keeping profit margins low in the short term to maximize future growth.

Interestingly, almost all of these relations appear to be a function of stable firm differ-

⁹ Note that prior to the firm's first open-source activity, it is classified as a non-open-source firm.

ences. When including firm fixed effects in Column (2), all variables (except annual returns) become statistically insignificant. A firm’s average characteristics therefore do not appear to significantly differ from before to after their first open-source activity. In combination with a relatively large adjusted R^2 for this regression, it seems that firm fixed effects can account for most of the differences between open-source and non-open-source firms.

In Columns (3) and (4), we focus only on open-source firms and examine the determinants of monthly open-source activity. To measure open-source activity, we use the number of commits to repositories owned by the firm in that month. Column (3) compares firms to their industry peers and reports similar results as those reported in Column (1). Specifically, firms with more open-source activities tend to be larger and more innovative, although they also tend to be less profitable and have lower returns. These firms also tend to have more market power, but receive fewer benefits from product-market network effects and face a more fluid product market.

However, the relations for product-market characteristics reverse when firm fixed effects are included in Column (4). For example, while firms with more market power tend to have more open-source activities, these firms are especially active when their market power is lower relative to their sample average. This result suggests that firms may use open-source activities to maintain their market power, however further research is required to make a stronger conclusion. We also find that firms engage in more open-source activities when their performance is relatively weak (Sales Growth, Annual Returns) and, intuitively, when they have more employees.

4 Open-Source Value

Having characterized open-source firms and assessed the determinants of open-source activity, we next turn to estimating the economic value of open-source innovation. Specifically, we leverage financial markets to estimate the value, in dollars, of GitHub repositories based

on stock returns around the date on which the repository was made public. This estimate corresponds to the private value (i.e., the value captured by the firm, as opposed to the value generated for all firms) of the repository as a whole (i.e., the sum of the value of the innovation plus the value of being open source). This estimate is also a forward-looking estimate of value as of the date that the repository was made public, and as such does not capture changes in value that may occur as the project is further developed. In the following, we outline the procedure to estimate repository value, validate these estimates using a realized measure of repository popularity and a placebo test, and investigate the determinants of open-source value.

4.1 Estimating value

Our procedure for estimating the value of repositories closely follows the procedure developed by [Kogan et al. \(2017\)](#) to estimate patent value and used by [Desai et al. \(2023\)](#) to estimate trademark value. We provide a detailed discussion of this procedure in [Internet Appendix C](#) and briefly outline the crucial points in this section.

The procedure involves observing stock returns in the three-day window following the announcement of the repository, $[t, t + 2]$. We cumulate market-adjusted returns over the three-day announcement window for repository i , which we label R_i . We assume that R_i is a function of both investor reaction to the repository announcement, v_i , and idiosyncratic noise.

We construct the estimate of repository value as the product of the investor reaction to the repository announcement and the firm’s market capitalization on the day prior to the announcement. If multiple repositories are announced on the same day, we assume the value is evenly distributed across those repositories. The value of repository i , ξ_i , is thus calculated as

$$\xi_i = \frac{1}{N_i} E[v_i | R_i] M_i, \tag{1}$$

where N_i is the number of repositories announced on that day, $E[v_i|R_i]$ is the expected return attributable to the repository announcement conditional on observing the three-day cumulative market-adjusted return R_i , and M_i is the market capitalization of the firm on the day prior to the repository announcement.

[Internet Appendix C](#) discusses our estimation of the conditional expected return in Equation (1), which adopts the same distributional assumptions as [Kogan et al. \(2017\)](#). Importantly, these assumptions imply that repositories have strictly positive values. While it is possible that open-source projects provide value to competitors that make the projects less valuable to the firm itself, we assume that firms will only choose to make projects open source if the net effect still results in a positive value for the firm.

4.2 Summary statistics

We estimate Equation (1) for the 29,543 original repositories that have available announcement dates as well as the required stock return data from CRSP.¹⁰ In Panel A of Table 3, we report the mean, standard deviation, and multiple distribution percentiles (1st, 5th, 10th, 25th, 50th, 75th, 90th, 95th, and 99th) of several variables. The mean (median) three-day cumulative market-adjusted announcement return (R_i) is 0.12% (0.04%) and the mean (median) expected return attributable to the repository announcement ($E[v_i|R_i]$) is 0.27% (0.14%). The difference between the mean and median for both variables indicates that market reactions to repository announcements are positively skewed.

We find that the mean value, ξ , for repositories in our sample is \$842,849, and the median value is \$562,022. There is also significant skewness in this variable, with the 99th percentile of repository value exceeding \$5,000,000 and the most valuable repositories exceeding \$12,000,000 (untabulated). These values are reported in 2023 dollars. For comparison, the mean value for patents granted over a similar period (2015-2023) is \$53 million in 2023

¹⁰ We focus on original (i.e., not forked) repositories because the release date for forked repositories is less clearly defined.

dollars, as calculated using data from [Kogan et al. \(2017\)](#). Given that patented innovation typically requires higher investment and benefits from the excludability provided by the patent system, we consider the relatively lower mean repository value to be plausible.

We also report statistics on several repository characteristics. First, we report the number of stars each repository has received as of February 2024, which we use to measure the realized popularity of a repository. Again, we observe significant skewness in the variable, with the mean repository being starred 212 times but the median repository being starred only 10 times.

Second, we report repository complementarity and novelty scores, which reflect how much the repository complements the firm’s commercial products and how novel the repository is compared to existing solutions, respectively. These scores are determined using ChatGPT to analyze repository information and are defined between zero and one. We find that the majority of repositories (54.1%) significantly complement the firm’s commercial products, with a complementarity score of at least 0.5. However, there is still a significant portion of standalone repositories, with 23.6% of repositories having a complementarity score of 0. Most repositories are also not rated as particularly novel, with the median repository having a novelty score of 0.3.

Third, we report the distribution of repository size, which measures the amount of data (code, images, etc.) in bytes contained in the repository and displays significant skewness. Finally, we report the cumulative number of issues opened for repositories as of December 31, 2023. Opened issues tend to convey suggestions from the community to improve the repository or fix errors, but popular repositories are also more likely to have issues opened in general. Thus, the skewness we observe in the number of issues opened likely comes from the same phenomenon as the skewness we observe in the number of stars.

Panel B of [Table 3](#) reports a summary of repository values based on firm industries. The majority of the repositories in our sample are created by firms in the “Computer Software” industry (59.8%) and the average value for repositories in this industry (\$781,835) is similar

to that of the complete sample. On average, repositories from the “Consumer Durables, Non-Durables, Wholesale, Retail, and Some Services” industry are the most valuable (\$1,090,190), make up the second largest group of repositories (25.2%), and have the highest rate of repositories having fully permissive licenses (88.7%). The “Healthcare, Medical Equipment, and Drugs” industry is notable for having the fewest repositories (97) and lowest rate of permissive repositories (41.2%) in our sample, again consistent with the importance of excludability for this industry.

Panel C of Table 3 reports the 10 firms with the most valuable repository portfolios as well as the total value of all repositories in our sample. Amazon.com Inc and Microsoft Corp have the most valuable repository portfolios, both nearing \$8 billion. For both firms, the majority of their repositories have fully permissive licenses. The remaining listed firms are also well-known technology firms with a focus on innovation, such as Alphabet Inc, Adobe Inc, and International Business Machines Corp. In total, we find that the repositories in our sample generated nearly \$25 billion of private value for public firms.

Panel D of Table 3 reports the 10 programming languages that generate the most value as classified by each repository’s main programming language. Python is the most commonly represented language (20.1% of repositories) and is associated with the most total repository value (\$6,853,266,877). Python also has the highest average and median repository value among the languages listed.^{11,12} Go and JavaScript are notable for having the largest skewness in repository value among the languages listed, and C++ and HTML are notable for having the smallest percentage of repositories with a fully permissive license among the languages listed.

We also explore how value varies across repository topics. As described in Section 3.1, we use ChatGPT to create topic scores, between zero and one, assessing the extent to which

¹¹ This statement includes Jupyter Notebook as a Python language because it is a web-based interface often used to work interactively with Python code.

¹² Other languages with a higher average repository value and at least 10 repositories include Cuda (\$2,327,971, 23 repositories), Bicep (\$1,468,167, 79 repositories), Swift (\$1,435,616, 359 repositories), and CMake (\$1,326,104, 25 repositories).

each repository relates to each topic. To estimate how much value a repository generates for each topic, we multiply each topic score by the repository value. Note that topic scores do not necessarily sum to one for a given repository, so these statistics should not be interpreted as a decomposition of repository value. Table 4 reports the mean, median, and total value generated by repositories with non-zero topic scores for each topic. For these repositories, the table also reports the mean topic score, number of repositories, and percent of repositories with fully permissive licenses.

We find that repositories in “Core AI and ML” and “AI Applications” are the most valuable on average. This result is in part due to repositories in these topics having high topic scores on average. However, even after adjusting for this (e.g., dividing Mean ξ by Mean Topic Score), these are still the most valuable topics. Repositories in “Software Engineering” and “Cloud Infrastructure and DevOps” generate the most total value, although this is attributable to the majority of repositories being at least partially associated with these topics. The “Advanced Data Analysis” topic is also notable for having a relatively high mean value for its mean topic score.

Finally, we investigate how the average repository value has evolved over time. In Figure 3, we plot the average ξ , in 2023 dollars, of repositories released each quarter from 2015 through 2023. Average repository values hovered around \$400,000 from 2015 through 2017 and then jumped to between \$600,000 and \$800,000 from 2018 through 2019. This increase coincides with Microsoft’s acquisition of GitHub, which was announced in June of 2018. Average repository values significantly increased again the first quarter of 2020, which is likely attributable to expectations of increased digitalization resulting from the COVID-19 shutdown. Since 2020, average repository values have hovered around \$1,000,000, and most recently peaked above \$1,300,000. Thus, the average repository value reported for the whole sample understates the value of open-source innovation in recent years.

4.3 Determinants of open-source value

We next turn to explore which characteristics most strongly correlate with open-source value. We investigate a broad set of repository, firm, and product market characteristics to give an extensive assessment of these correlations. Many of these characteristics overlap with each other, so we also include them together in regressions to assess their marginal correlations with open-source value. [Internet Appendix D](#) reports univariate correlations among all pairs of variables included in our analysis. We view our results in this section as descriptive in nature and intended to provide a more complete characterization of open-source value.

4.3.1 Repository popularity

We begin by investigating repository popularity. While we contend that the statistics reported in [Table 3](#) indicate our measure of repository value is reasonable, there is still the possibility that significant noise in the estimation procedure renders the estimates largely uninformative. If this is the case, then we would expect repository value to be at most weakly correlated with subsequent repository popularity. We therefore investigate this possibility to provide validity for the estimates of repository value.

To measure repository popularity, we use the number of stars each repository has received as of February 2024. “Starring” a repository bookmarks it for the external user, which allows the user to stay updated on any changes made to the repository and indicates significant interest in the repository. We regress the natural logarithm of repository value, ξ , on the natural logarithm of one plus the number of stars the repository has received.¹³ We control for the natural logarithms of market capitalization (measured as of the repository announcement), volatility (measured over the announcement year), employees, and patent-portfolio value (both measured as of the prior year). We also include various fixed effects depending

¹³ We take the natural logarithm of each variable to adjust for the skewness documented in the previous section.

on the specification, including year, industry (at the three-digit SIC level), industry-year, firm, and firm-year fixed effects. We double-cluster standard errors by year and industry and all independent variables are standardized to facilitate interpretation.

The results of these regressions are reported in Table 5. In Column (1), we report the regression with year fixed effects and find a significant correlation between repository value and subsequent popularity. The regression reported in Column (2) adds industry fixed effects and the regression reported in Column (3) replaces these fixed effects with industry-year fixed effects. In both cases, we continue to find a significant correlation between repository value and subsequent popularity, with the correlation generally increasing in significance as the fixed effects become stricter. Economically, the estimate reported in Column (3) indicates that repositories that end up being one standard deviation more popular have an 8.9% higher valuation when released. Finally, Column (4) reports the regression with firm fixed effects, and Column (5) includes firm-year fixed effects and clusters standard errors by firm.¹⁴ Both results demonstrate a strong correlation between within-firm repository value and subsequent popularity.

We therefore conclude that repositories that are estimated to be more valuable when they are released tend to be significantly more popular in the future. This result indicates that the estimation procedure captures value-relevant information.

4.3.2 Placebo test

To provide further validation for the estimates of repository value, we perform a placebo test on the repository release date. While the strong correlation between repository value and future popularity suggests that investors pay attention to the releases of repositories on GitHub, this may be only a spurious relation.

We investigate this possibility by randomly assigning each repository a placebo release

¹⁴ We use firm clustering with firm-year fixed effects to match the placebo test procedure discussed in Section 4.3.2. We find similar results when double clustering standard errors by industry and year with firm-year fixed effects.

date in the same year as the true release date. We then estimate the placebo value of the repository using the market reaction on this placebo release date. Finally, we regress repository placebo value on the true number of stars subsequently received by the repository. The regression specification follows that of Column (5) of Table 5 (i.e., including firm-year fixed effects). We repeat this process 500 times. The resulting distribution of coefficient estimates, corresponding to number of stars, and t-statistics are plotted in Panel A and Panel B of Figure 4, respectively. Each panel also plots a vertical dotted line corresponding to the results from the identical regression with the true repository values (i.e., Column (5) of Table 5).

It is readily apparent in the figure that the true coefficient and t-statistic are outliers relative to the placebo estimates. While four of the 500 iterations produce t-statistics of similar magnitude to the true relation, none of the iterations produce a coefficient that approaches that of the true relation. We therefore conclude that the market reactions on repository release days are in fact, on average, directly related to the repository release.

4.3.3 Repository, firm, and product market characteristics

We next investigate the correlations between repository value and other repository characteristics, firm characteristics, and product market characteristics. The results are reported in Table 6 with Panels A, B, and C corresponding to each category of characteristics, respectively. Each regression includes industry-year fixed effects and controls for repository popularity, stock market capitalization, stock volatility, employees, and total patent value. Within each panel, characteristics are sequentially introduced, with the final column reporting the regression including all characteristics from that category.

In Panel A, we investigate repository characteristics. We first consider the type of license covering the repository. Each repository is classified into one of three groups based on its license: permissive (no restrictions on use), copyleft (some restrictions on use), and other (cannot be classified). We then include indicator variables for “copyleft” repositories and

“other” repositories in the regression such that “permissive” repositories represent the omitted category. Traditional models of the economics of innovation suggest that excludability increases the private value of innovation (Schumpeter, 1912), and Lerner and Tirole (2005b) explicitly discuss how license restrictiveness increases open-source value. Consistent with these theories, we find that copyleft repositories are approximately 10.5% more valuable, on average, than permissive repositories (Column (8)).

We next examine repositories designated as templates. Template repositories can be easily duplicated into new repositories with identical directory structure, branches, and files without keeping the commit history. Given that these repositories are less likely to contain a unique innovation, we expect them to be less valuable on average. Consistent with this notion, we find that template repositories are approximately 17.5% less valuable, on average, than non-template repositories (Column (8)).

We also examine how repository value correlates with repository complementarity and novelty. On the one hand, repositories that complement the firm’s commercial products may be more valuable because they drive demand for those commercial products, thus increasing profitability for the firm. On the other hand, standalone (i.e., less complementary) repositories may be more valuable because they represent a more substantial project undertaken by the firm. Consistent with the latter interpretation, we find that complementarity is negatively associated with repository value. However, it is possible that complementarity also (inversely) captures an aspect of novelty, which we expect to be positively associated with repository value. Consistent with this prediction, our measure of novelty is positively related to repository value in Columns (4) and (8). Specifically, a typical repository is 38.3% more valuable when it has a novelty score of 1.0 compared to 0.0. Moreover, when controlling for novelty in Column (8), we find that complementarity remains negatively related to repository value: a typical repository is 17.5% less valuable when it has a complementarity score of 1.0 compared to 0.0. Thus, we conclude that complementarity is not, in the cross section,

a first-order driver of repository value.¹⁵

We next examine the byte size of repositories. Repositories with more lines of code will have a larger byte size, all else equal, and may therefore be more valuable. However, we find a negative relation between repository size and repository value that is statistically significant when included by itself (Column (5)) and statistically insignificant when controlling for other repository characteristics (Column (8)). This result could be due to other files included in the repository, such as images, that increase the repository size and do not represent additional lines of code. To investigate this possibility, we examine the subset of repositories (2,223 repositories) for which byte size is separately categorized into binary (e.g., images) and non-binary (e.g., lines of code) data. However, we still find that both types of data are negatively and insignificantly related to repository value (untabulated).¹⁶ We therefore conclude that larger repositories are not necessarily more valuable.

We then examine the number of repositories previously released by the firm. This quantity is negatively and significantly related to repository value, suggesting that firms producing a lower quantity of repositories tend to produce higher-quality repositories. Finally, we examine the cumulative number of issues opened for the repository as of December 31, 2023. While this value captures potential bugs or problems with the repository, it also scales with the overall popularity of the repository. However, we control for repository popularity in the regression using the number of stars, and thus interpret issues opened as a negative reflection of repository quality. Consistent with this interpretation, we find a negative and significant

¹⁵ We cannot exclude the possibility that spillover profitability from complementarity incentivizes firms to make highly complementary repositories open source as opposed to closed source. Thus, complementarity may still be a driver of open-source value, but only for repositories with high complementarity, which tend to be relatively less valuable. It is also possible that a firm’s portfolio of repositories is made up of many high-complementarity repositories and relatively few standalone repositories, such that complementarity provides significant value across the firm’s whole portfolio. However, we find that repositories with a complementarity score of at least 0.5 represent only 27.9%, on average, of the total value of firms’ portfolios of repositories, and only 7.7% for the median firm. It therefore does not appear that complementarity represents a significant portion of the value of firms’ repository portfolios.

¹⁶ We do, however, find that the ratio of non-binary byte size to total byte size is positively and significantly related to repository value (untabulated). It therefore appears that non-binary data (e.g., lines of code) contribute more to the value of a repository than binary data (e.g., images).

relation between the number of issues opened in the future and repository value estimated by investors at announcement.

In Panel B of Table 6, we investigate firm characteristics. These characteristics include market-to-book ratio, return-on-assets, investment, annual stock return, annual sales growth, tangibility, and R&D expenditure. For observations with missing R&D expenditure, we replace the value with zero and set an indicator variable, R&D Exp Missing, equal to one. We find that the majority of these variables are not significantly related to repository value in the cross section. The two exceptions are return-on-assets, which suggests that more-profitable firms tend to produce more-valuable repositories, and the indicator variable for missing R&D expenditure, which suggests that such firms also tend to produce more-valuable repositories. The latter result could imply that investment in open-source innovation is difficult for firms to quantify as R&D expenditure. At the least, the result demonstrates that R&D expenditure does not fully capture open-source activities. More generally, the lack of statistical significance across firm characteristics after including our standard controls (market capitalization, volatility, number of employees, and patent value), particularly in contrast to the results for repository characteristics, suggests these controls capture much of the firm-level variation in repository value. This finding narrows the scope for potentially omitted variables that could confound our analysis of firm growth in the next section.

Finally, in Panel C of Table 6, we investigate product market characteristics. These variables are of particular interest because a firm’s product-market rivals are most likely to benefit from the open-source nature of a firm’s repositories. With this in mind, we distinguish between the private value of a repository, which is captured by the innovating firm, and the public value of the repository, which is captured by other firms. The total value of the repository is thus the sum of the private and public values.

We first examine market power, as measured by the estimate of markups developed by [Pellegrino \(2024\)](#). Market power is negatively related to competition in theory and reflects a firm’s ability to extract rents from its customers. Firms with more market power

may face less risk of their repositories being used by competitors, which could either make them more willing to share more-valuable innovation via open source or allow them to extract more private value from the repository. Consistent with these two possibilities, we find that firms with more market power tend to produce more-valuable repositories.

We then consider the firm's centrality in the product market network. To the extent that this quantity reflects the level of competition faced by the firm, similar to market power, we would expect centrality to be negatively related to repository value. However, centrality also measures the extent to which the firm benefits from network effects in the product market, which we expect to be positively related to repository value. Consistent with these opposing predictions, we find that centrality is insignificantly related to repository value when included by itself in Column (2). However, this relation becomes positive and significant in Column (6) when including other product market characteristics that also reflect competition (e.g., market power). The result in Column (6) therefore supports the hypothesis that network effects are a significant driver of value for open-source innovation.

Finally, we also consider the scope, product market similarity, and product market fluidity of firms. We find that scope is negatively related to repository value, suggesting that firms with a sharper product focus tend to produce more-valuable repositories. We also find that product market fluidity is negatively related to repository value, suggesting that repositories create more private value when the firm is in a more stable product market.

In summary, the results for product market characteristics suggest that repository value is negatively related to competition. Given that our estimates of value reflect private value, it is possible that firms facing less competition are able to capture a larger fraction of the total value created by the repository. Alternatively, firms facing less competition may be more willing to share innovation that has a higher total value. In either case, we find that competition is an important consideration for firms producing open-source innovation.

Finally, it is important to note that repository popularity is positively and significantly correlated with repository value across all regressions reported in Table 6. This result further

supports the validation exercises from the previous sections.

5 Open-Source Innovation and Firm Growth

Technological innovations are recognized for driving long-term growth for companies, usually in terms of the quality or number of products (e.g., [Aghion and Howitt \(1992\)](#)). However, the financial gains associated with that growth depend on the excludability associated with the innovation. Given that open-source licenses grant usage rights to the general public, the relation between open-source innovation and firm growth, in financial terms, is uncertain.

To investigate this relation, we first calculate the firm-level repository value, $\xi_{f,t}$, as the sum of the repository values ξ_i for all repositories posted by firm f in year t . We consider the growth of several dependent variables (Y) including sales, profits, the number of employees, and the value and number of patents:

$$\ln Y_{f,t+k} - \ln Y_{f,t} = \beta_k \ln(\xi_{f,t} + 1) + \psi_k X_{f,t} + \epsilon_{f,t+k}, \quad (2)$$

where the horizon k varies from one to three years. The vector X includes the natural logarithms of market capitalization, idiosyncratic volatility, number of employees, patent-portfolio value, and $Y_{f,t}$. We also include industry (at the three-digit SIC level) and year fixed effects and double cluster standard errors by industry and year.

The results are outlined in [Table 7](#). In this analysis, our focus is on the coefficient β_k , which captures the impact of open-source innovations by the firm on firm growth. The three columns present our estimates of β_k for $k = 1, 2$, and 3 , respectively. We observe a significant positive relation between a company's open-source innovation and its future growth in terms of sales, profits, and number of employees, over the following three years. Economically, the coefficients imply that over the following three years, a 100% greater firm-level repository value stimulates the growth of sales by 80 basis points, the growth of

profits by 70 basis points, and the growth of employment by 80 basis points. We also observe complementarity between open-source innovation and patentable innovation. A 100% greater firm-level repository value is associated with a 2.1 percentage point increase in the value of new patents granted and a 1.6 percentage point increase in the number of new patents granted over the following three years. Furthermore, the larger coefficients for patent value relative to the number of patents indicate that the average value per patent increases as well. Overall, our findings suggest that even when innovation is made accessible to others, firms experience benefits from that innovation.

6 Conclusion

Given the importance of excludability in generating private value from innovation, the growing involvement of public firms in open-source innovation is initially surprising. To explore this puzzling phenomenon, we construct an extensive dataset of open-source activities by public firms on GitHub, the largest open-source development platform, and use financial markets to develop a measure of the private value of open-source innovation.

We find that open-source engagement is highly prevalent in the U.S. economy. Despite only 18% of public firms having open-source projects, those firms represent 68% of stock-market value and 80% of R&D expenditure across 86% of Fama-French 49 industries. Firms with open-source projects tend to be larger, more valuable, more innovative, and face less competition on average.

We estimate the private value of open-source projects based on stock-market reactions. The average project in our sample is valued at \$842,849 (in 2023 dollars), with average values exceeding \$1,300,000 by the end of 2023. The total private value created by all projects in our sample is \$25 billion. We find that projects with licenses that provide at least some restriction in use generate more private value, which highlights that excludability is still important in the open-source setting. We also find that standalone projects are

more valuable than projects that complement the firm's commercial products. This result suggests that at least in the cross section of open-source projects, the potential for freely available open-source projects to promote the adoption of commercial products is not a first-order driver of private value. Furthermore, larger projects (e.g., more lines of code) do not necessarily create more private value and firms facing less competition tend to generate more private value from their open-source projects.

Finally, we find that valuable open-source innovation predicts future firm growth in terms of sales, profits, number of employees, and both the number and value of new patents. Thus, despite the innovation being made available to others, firms still benefit from open-source innovation.

In summary, these results provide new evidence of the value of innovation without excludability. Our estimates of open-source value open up avenues for future research on the benefits firms gain from open-source innovation and the valuation of intangible assets.

References

- Aghion, P. and P. Howitt (1992). A Model of Growth Through Creative Destruction. *Econometrica* 60(2), 323–351.
- Ahmadi, A., A. Kecskés, R. Michaely, and P.-A. Nguyen (2024). Producing AI Innovation and Its Value Implications. Working paper, York University, University of Hong Kong.
- Alexy, O., J. West, H. Klapper, and M. Reitzig (2018). Surrendering control to gain advantage: Reconciling openness and the resource-based view of the firm. *Strategic Management Journal* 39(6), 1704–1727.
- Allen, R. C. (1983). Collective invention. *Journal of Economic Behavior & Organization* 4(1), 1–24.
- Arrow, K. (1962). Economic Welfare and the Allocation of Resources for Invention. In *The Rate and Direction of Inventive Activity: Economic and Social Factors*, pp. 609–626. Princeton University Press.
- Austin, D. H. (1993). An Event-Study Approach to Measuring Innovative Output: The Case of Biotechnology. *American Economic Review* 83(2), 253–258.
- Blind, K., M. Böhm, P. Grzegorzewska, A. Katz, S. Muto, S. Pättsch, and T. Schubert (2021). The impact of Open Source Software and Hardware on technological independence, competitiveness and innovation in the EU economy. European Commission, Ed.
- Chen, M. A., Q. Wu, and B. Yang (2019, May). How Valuable Is FinTech Innovation? *Review of Financial Studies* 32(5), 2062–2106.
- Conti, A., C. Peukert, and M. Roche (2021). Beefing IT up for your Investor? Open Sourcing and Startup Funding: Evidence from GitHub. Accepted at *Organization Science*.
- Crouzet, N., J. C. Eberly, A. L. Eisfeldt, and D. Papanikolaou (2022). The Economics of Intangible Capital. *Journal of Economic Perspectives* 36(3), 29–52.
- Dahlander, L. and D. M. Gann (2010). How open is innovation? *Research Policy* 39(6), 699–709.
- Dahlander, L., D. M. Gann, and M. W. Wallin (2021). How open is innovation? A retrospective and ideas forward. *Research Policy* 50(4), 104218.
- Davis, J. L., E. F. Fama, and K. R. French (2000). Characteristics, Covariances, and Average Returns: 1929 to 1997. *Journal of Finance* 55(1), 389–406.
- Desai, P., E. Gavrilova, R. Silva, and M. Soares (2023). The Value of Trademarks. Working Paper, Nova School of Business and Economics.
- Fama, E. F. and K. R. French (1997). Industry costs of equity. *Journal of Financial Economics* 43(2), 153–193.

- Goldfarb, A. and C. Tucker (2019). Digital Economics. *Journal of Economic Literature* 57(1), 3–43.
- Gómez-Cram, R. and A. Lawrence (2024). The Value of Software. Working paper.
- Greenstein, S. and F. Nagle (2014). Digital dark matter and the economic contribution of Apache. *Research Policy* 43(4), 623–631.
- Hall, B. H., A. Jaffe, and M. Trajtenberg (2005). Market Value and Patent Citations. *RAND Journal of Economics* 36(1), 16–38.
- Harhoff, D., J. Henkel, and E. von Hippel (2003). Profiting from voluntary information spillovers: how users benefit by freely revealing their innovations. *Research Policy* 32(10), 1753–1769.
- Henkel, J., S. Schöberl, and O. Alexy (2014). The emergence of openness: How and why firms adopt selective revealing in open innovation. *Research Policy* 43(5), 879–890.
- Hoberg, G. and G. Phillips (2016). Text-Based Network Industries and Endogenous Product Differentiation. *Journal of Political Economy* 124(5), 1423–1465.
- Hoberg, G. and G. Phillips (2024). Scope, Scale and Concentration: The 21st Century Firm. *Journal of Finance*. Forthcoming.
- Hoberg, G., G. Phillips, and N. Prabhala (2014). Product Market Threats, Payouts, and Financial Flexibility. *Journal of Finance* 69(1), 293–324.
- Hoffmann, M., F. Nagle, and Y. Zhou (2024). The value of open source software. Working paper, Harvard University, University of Toronto.
- Kogan, L., D. Papanikolaou, A. Seru, and N. Stoffman (2017). Technological Innovation, Resource Allocation, and Growth. *Quarterly Journal of Economics* 132(2), 665–712.
- Lerner, J. and J. Tirole (2002). Some Simple Economics of Open Source. *Journal of Industrial Economics* 50(2), 197–234.
- Lerner, J. and J. Tirole (2005a). The Economics of Technology Sharing: Open Source and Beyond. *Journal of Economic Perspectives* 19(2), 99–120.
- Lerner, J. and J. Tirole (2005b). The Scope of Open Source Licensing. *Journal of Law, Economics, & Organization* 21(1), 20–56.
- Lin, Y.-K. and L. M. Maruping (2022). Open Source Collaboration in Digital Entrepreneurship. *Organization Science* 33(1), 212–230.
- Murciano-Goroff, R., R. Zhuo, and S. Greenstein (2021). Hidden software and veiled value creation: Illustrations from server software usage. *Research Policy* 50(9), 104333.
- Nagle, F. (2018). Learning by Contributing: Gaining Competitive Advantage Through Contribution to Crowdsourced Public Goods. *Organization Science* 29(4), 569–587.

- Nagle, F. (2019). Open Source Software and Firm Productivity. *Management Science* 65(3), 1191–1215.
- Pakes, A. (1985). On Patents, R&D, and the Stock Market Rate of Return. *Journal of Political Economy* 93(2), 390–409.
- Parker, G., M. Van Alstyne, and X. Jiang (2017). Platform Ecosystems: How Developers Invert the Firm. *MIS Quarterly* 41(1), 255–266.
- Pellegrino, B. (2024). Product Differentiation and Oligopoly: a Network Approach. *American Economic Review*. Forthcoming.
- Robbins, C., G. Korkmaz, L. Guci, J. B. S. Calderón, and B. Kramer (2021). A First Look at Open-Source Software Investment in the United States and in Other Countries, 2009-2019. Paper prepared for the IARIW-ESCoE Conference.
- Sas, C., A. Capiluppi, C. Di Sipio, J. Di Rocco, and D. Di Ruscio (2023). Gitranking: A ranking of github topics for software classification using active sampling. *Software: Practice and Experience* 53(10), 1982–2006.
- Schumpeter, J. (1912). *The Theory of Economic Development*. Cambridge, MA: Harvard University Press.
- Teece, D. J. (2018). Profiting from innovation in the digital economy: Enabling technologies, standards, and licensing models in the wireless world. *Research Policy* 47(8), 1367–1387.
- von Hippel, E. and G. von Krogh (2003). Open Source Software and the “Private-Collective” Innovation Model: Issues for Organization Science. *Organization Science* 14(2), 107–225.

Figure 1
Trends in Open-Source Engagement Among U.S. Public Firms (2015-2023)

This graph plots the time series of U.S. firms' participation in open-source activities through the creation of public GitHub repositories from 2015 to 2023. It represents the proportion of firms making repositories public in terms of total number of firms, market capitalization, and R&D expenditure (left y-axis), and it tracks the cumulative number of public repositories owned by these firms (right y-axis).

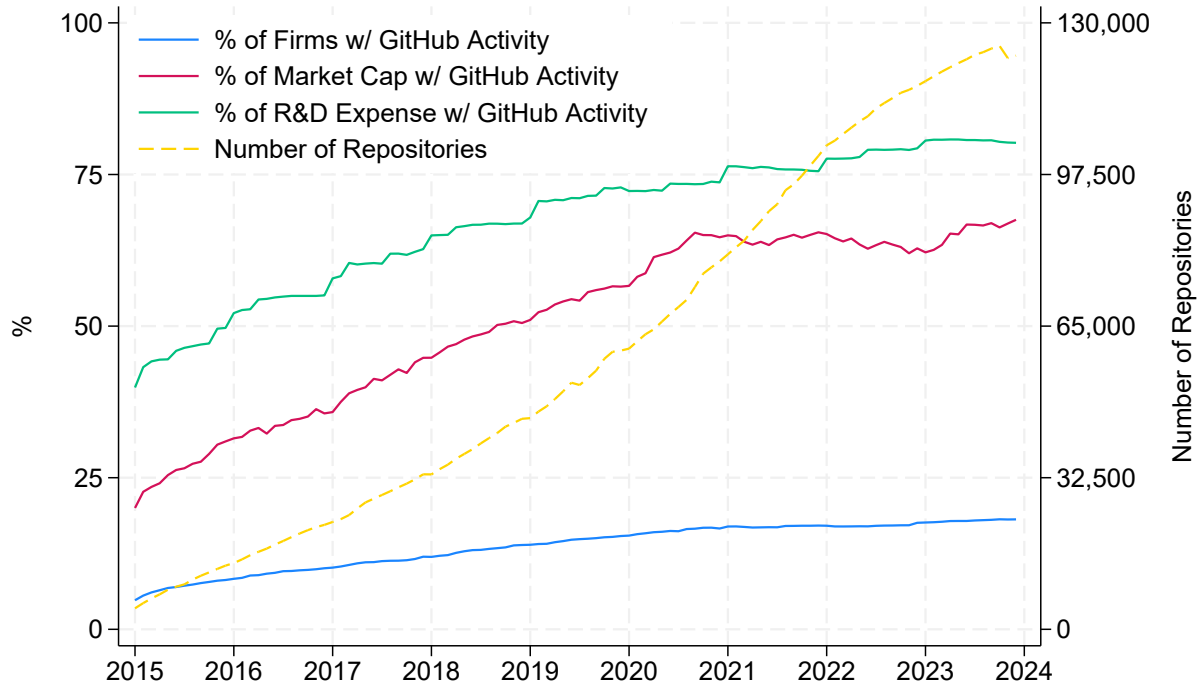


Figure 2
Industry Distribution of Open-Source Engagement

This figure illustrates the percentage of firms with GitHub activity and their respective repositories across various industries over the period from 2015 to 2023. The percentages are derived by dividing the number of firms with GitHub activity in a given industry by the total number of firms active on GitHub. The left pie chart shows the proportion of firms with any GitHub activity, categorized by industry, while the right pie chart displays the distribution of all repositories owned by these firms. Industry classification adheres to the Fama-French 5 Industries, with the exception that the Computer Software and Finance industries are distinctly separated using the Fama-French 49 Industries.

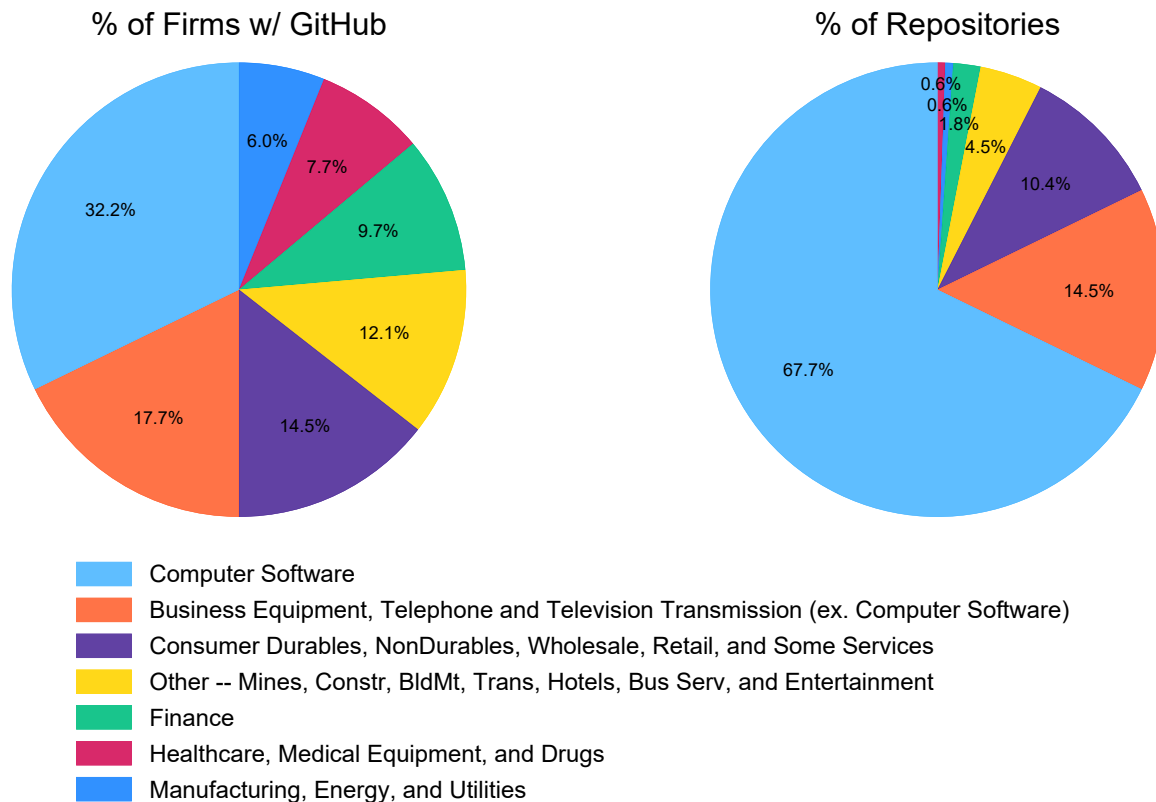


Figure 3
Average Estimated Repository Value by Quarter

This figure displays the trajectory of the average repository value, ξ , from 2015 to 2023. The values are computed using the methodology detailed in Section 4.1 and have been adjusted to reflect 2023 dollar values.

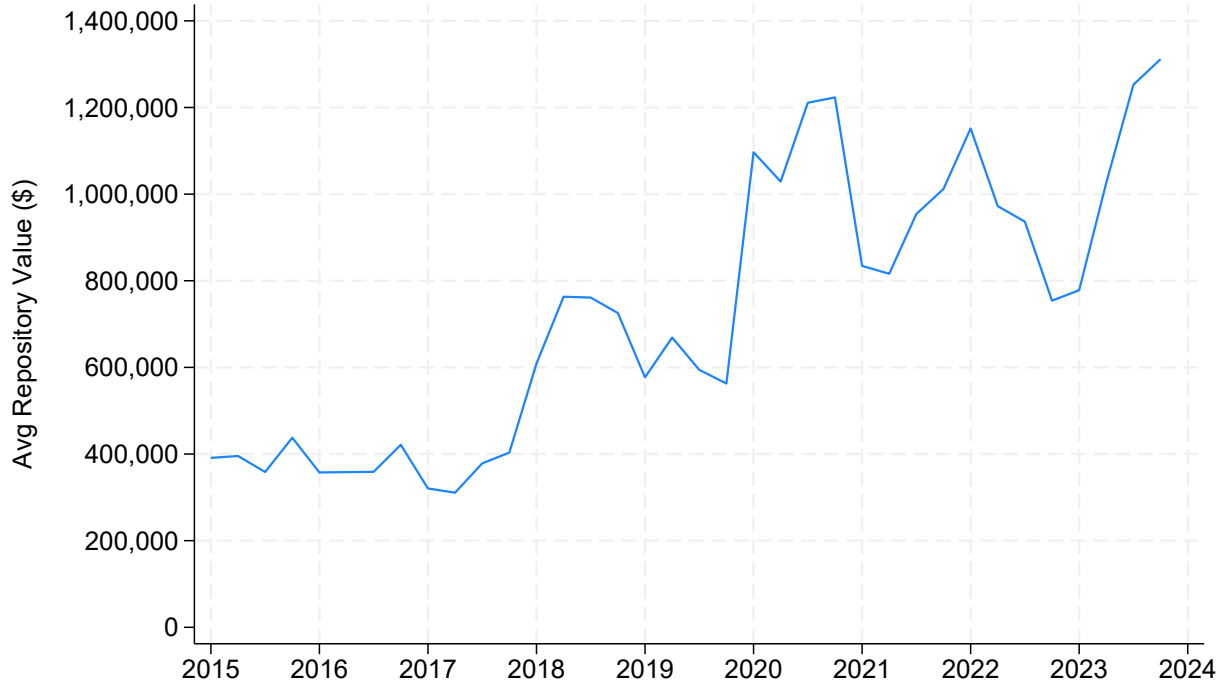


Figure 4
Placebo Test Results

This figure displays the distributions of coefficient estimates and t-statistics from 500 iterations of placebo tests, conducted to validate the measure of private value for repositories. In the procedure, each repository is assigned a random placebo release date within its actual release year. The repository's placebo value is then determined by the market response on that date. The upper figure shows the distribution of coefficient estimates linked to the number of stars received, while the lower figure shows the distribution of corresponding t-statistics. Vertical dotted lines in both panels mark the actual coefficient estimate and t-statistics for comparison.

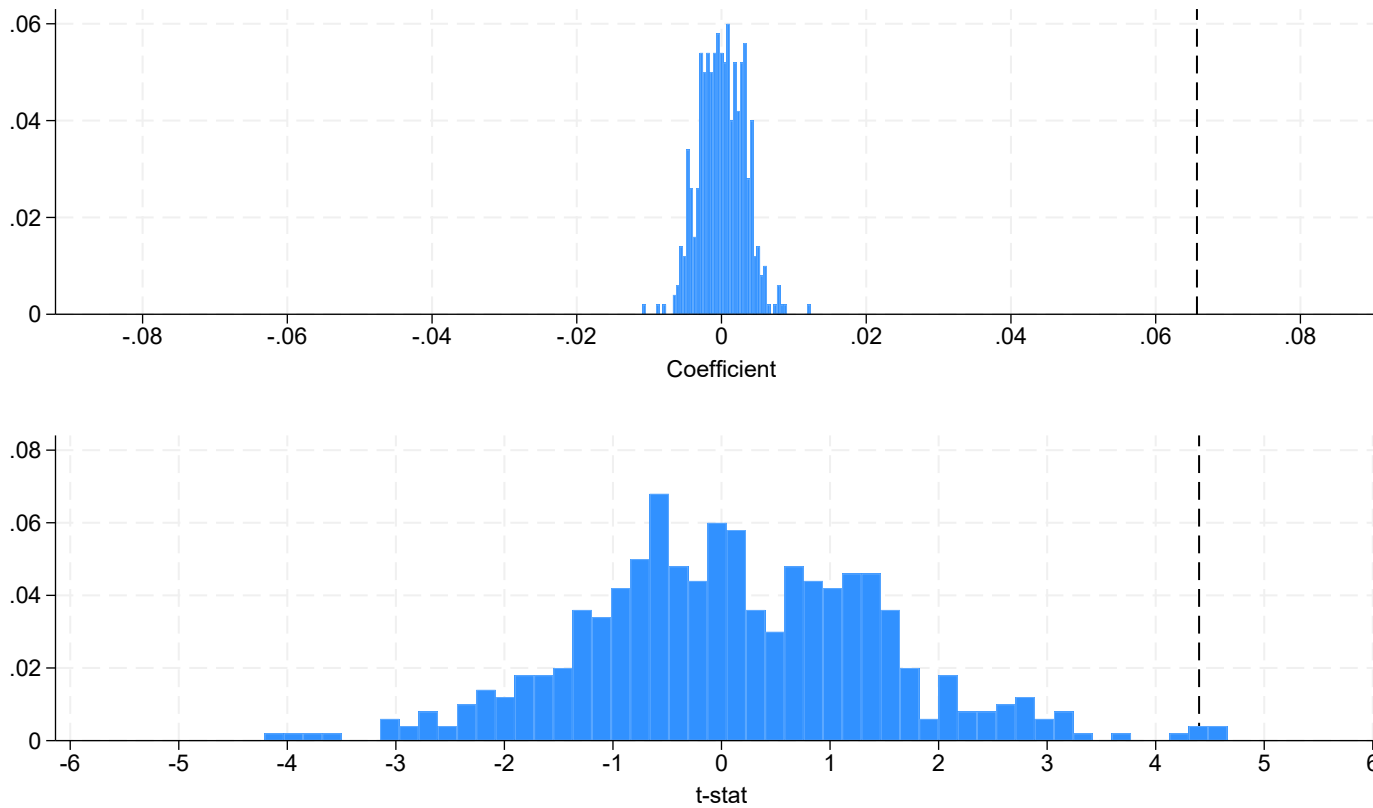


Table 1
Summary of GitHub Activity

Panel A presents the prevalence of GitHub activity among firms, segmented by industry based on the modified Fama-French 5 Industries classification, with an explicit distinction made for the Computer Software and Finance industries. Panel A details the percentage of firms engaged in GitHub activities and the distribution of repository ownership within each industry. Panel B compares key financial characteristics from 2015 to 2023 between firms that are active on GitHub and firms that are not. See Table A1 for the definition of variables.

Panel A: GitHub Activity

	% GitHub	Number of Repositories								Total	N Firms
		Mean	Std	p25	p50	p75	p90	p95	p99		
Total	18.1%	30.9	430.1	0	0	0	17	62	425	122,971	3,982
Consumer Durables, NonDurables, Wholesale, Retail, and Some Services	19.0%	23.6	393.4	0	0	0	13	36	197	12,757	541
Manufacturing, Energy, and Utilities	7.5%	1.2	6.4	0	0	0	0	6	38	6961	575
Business Equipment, Telephone and Television Transmission (ex. Computer Software)	34.9%	49.1	229.6	0	0	8	87	156	1253	17,726	361
Healthcare, Medical Equipment, and Drugs	6.4%	0.8	5.7	0	0	0	0	2	24	720	861
Other – Mines, Constr, BldMt, Trans, Hotels, Bus Serv, and Entertainment	20.6%	13.3	65.8	0	0	0	14	46	336	5,558	418
Computer Software	62.6%	226.5	1298.7	0	13	82	340	657	4551	82,896	366
Finance	10.8%	3.4	17.3	0	0	0	2	18	93	2,173	636

Panel B: Financial Statistics

	GitHub Firms		Non-GitHub Firms	
	Mean	Median	Mean	Median
Market Capitalization	32,220,913	3,743,408	1,418,477	71,523
Employees	32.5	4.7	7.9	1.1
Number of Patents	1,264	13	67	0
Market-to-Book	6.26	3.49	2.69	1.55
Return-on-Assets	-1.39%	2.26%	-1.35%	3.25%
Investment	3.34%	2.06%	7.06%	4.30%
Annual Returns	14.41%	6.50%	15.21%	5.74%
Sales Growth	14.93%	9.11%	17.27%	8.88%
Tangibility	14.87%	9.02%	27.35%	20.86%
R&D Exp / Total Assets	8.22%	4.79%	3.99%	0.00%
Market Power	3.13	2.21	2.29	1.64
Scope	11	10	8	7
Product Market Centrality	0.0043	0.0024	0.0086	0.0039
Product Market Similarity	4.16	1.74	11.64	2.00
Product Market Fluidity	5.23	4.91	7.67	6.92

Table 2
Determinants of Open-Source Activity

This table reports regression results to examine the factors influencing the extensive and intensive margins of GitHub activity among U.S. public firms. The dependent variable in Columns (1)-(2), *Github*, is a dummy variable that equals one from the first month a firm engages in any open-source activity (open-source firm) on GitHub. In Columns (3)-(4), the dependent variable is the natural logarithm of one plus the number of commits made to repositories owned by the open-source firm each month. See Table A1 for the definition of variables. Standard errors double clustered by industry and year are reported in parentheses. ***, **, and * indicate significance at the 1, 5, and 10% levels, respectively.

	(1) GitHub	(2) GitHub	(3) ln(N Commits + 1)	(4) ln(N Commits + 1)
ln(Mkt Cap)	0.087*** (0.018)	0.016* (0.009)	1.127*** (0.272)	0.026 (0.223)
ln(Employees)	0.003 (0.011)	0.024 (0.019)	-0.092 (0.248)	0.532* (0.273)
ln(N Patents + 1)	0.043*** (0.012)	0.017 (0.042)	0.613*** (0.138)	0.035 (0.249)
Market-to-Book	0.014* (0.008)	0.002 (0.003)	-0.025 (0.050)	0.081* (0.047)
Return-on-Assets	-0.008* (0.004)	0.002 (0.002)	-0.178*** (0.060)	0.017 (0.049)
Investment	0.002 (0.004)	-0.000 (0.002)	0.014 (0.128)	0.101 (0.101)
Return (t-12 to t-1)	-0.007*** (0.002)	-0.003*** (0.001)	-0.072** (0.031)	-0.041** (0.017)
Sales Growth	-0.001 (0.004)	-0.001 (0.002)	0.065 (0.042)	-0.060* (0.036)
Tangibility	-0.003 (0.010)	-0.010 (0.007)	0.130 (0.174)	-0.015 (0.178)
R&D Exp/Total Assets	0.029** (0.012)	-0.003 (0.005)	0.280* (0.149)	0.056 (0.091)
R&D Exp Missing	-0.041** (0.018)	-0.014 (0.016)	0.416* (0.234)	-0.777*** (0.218)
Market Power	0.009 (0.007)	0.001 (0.004)	0.215*** (0.044)	-0.117*** (0.030)
Scope	-0.005 (0.010)	0.003 (0.005)	0.066 (0.122)	-0.050 (0.085)
Product Market Centrality	-0.037** (0.016)	-0.001 (0.015)	-0.472** (0.213)	0.168* (0.092)
Product Market Similarity	0.002 (0.017)	-0.005 (0.011)	0.149 (0.239)	-0.078** (0.034)
Product Market Fluidity	0.009 (0.007)	-0.006* (0.003)	0.120** (0.052)	-0.079* (0.046)
Observations	208,528	208,513	26,422	26,413
Adj. R ²	0.331	0.866	0.327	0.781
Industry x Time FE	✓	✓	✓	✓
Firm FE		✓		✓
Sample	All firms	All firms	GitHub = 1	GitHub = 1

Table 3
Summary of Repository Value

This table reports summary statistics of repository values estimated through the methodology detailed in Section 4.1 over the period from 2015 to 2023. Panel A summarizes announcement returns, repository values, and other repository characteristics. R is the three-day cumulative market-adjusted announcement return. $E[v|R]$ is the conditional expected return attributable to the repository announcement. ξ is the estimated repository value reported in 2023 dollars. See Table A1 for definitions of the remaining variables. Panel B provides a breakdown of repository values by industry, adhering to a modified version of the Fama-French 5 Industries classification and distinguishing the Computer Software and Finance industries separately. % Permissive is the percent of repositories with permissive licenses. Panel C lists the top 10 firms based on the aggregate value of their GitHub repository portfolios. Panel D lists the top 10 programming languages based on their aggregate value.

Panel A: Summary Statistics

	Mean	Std	p1	p5	p10	p25	p50	p75	p90	p95	p99	N
R	0.12%	3.49%	-9.17%	-4.59%	-3.06%	-1.31%	0.04%	1.48%	3.47%	5.03%	10.13%	29,543
$E[r R]$	0.27%	0.31%	0.01%	0.03%	0.04%	0.06%	0.14%	0.39%	0.68%	0.89%	1.36%	29,543
ξ	842,849	1,044,529	1,637	9,514	23,006	140,577	562,022	1,146,063	1,953,618	2,742,288	5,161,436	29,543
Stars	212.2	2,226.6	0	0	0	2	10	44	207	565	3752	29,543
Complementarity	0.41	0.29	0	0	0	0.1	0.5	0.6	0.8	0.8	0.8	29,508
Novelty	0.26	0.14	0	0.1	0.1	0.2	0.3	0.3	0.5	0.5	0.6	29,529
Repo Size	31,322.1	282,597.8	5	14	27	113	808	6374	39712	100330	526113	29,535
N Issues Opened	56.0	1,098.4	0	0	0	0	1	8	43	123	814	29,543

Panel B: ξ by Industry

	Total ξ	Mean ξ	Median ξ	N Repos	% Permissive
Computer Software	13,472,579,594	781,835	463,577	17,232	68.3%
Consumer Durables, NonDurables, Wholesale, Retail, and Some Services	7,911,507,157	1,090,190	892,361	7,257	88.7%
Business Equipment, Telephone and Television Transmission (ex. Computer Software)	3,091,412,879	961,559	232,003	3,215	49.4%
Other – Mines, Constr, BldMt, Trans, Hotels, Bus Serv, and Entertainment	181,271,280	345,938	216,893	524	65.3%
Finance	89,538,721	255,825	140,748	350	77.7%
Healthcare, Medical Equipment, and Drugs	51,579,235	531,745	269,625	97	41.2%
Manufacturing, Energy, and Utilities	44,608,886	273,674	213,635	163	62.0%

Panel C: Firms with Most Valuable GitHub Portfolios

	Portfolio ξ	Mean ξ	Median ξ	N Repos	% Permissive
Amazon.com Inc	7,814,250,435	1,145,281	924,267	6,823	91.5%
Microsoft Corp	7,759,843,486	1,129,855	862,161	6,868	72.7%
Meta Platforms Inc	2,280,475,890	1,949,125	1,738,289	1,170	45.0%
Alphabet Inc	1,763,699,392	1,031,403	808,095	1,710	86.8%
NVIDIA Corporation	1,391,421,878	2,394,874	1,868,125	581	49.6%
Apple Inc	1,059,601,598	4,489,837	3,702,204	236	55.1%
Salesforce Inc	495,777,701	477,168	361,598	1,039	74.4%
Adobe Inc	225,007,179	646,572	578,176	348	77.6%
International Business Machines Corp	211,208,262	344,549	334,237	613	51.7%
Oracle Corp	206,261,704	661,095	612,226	312	71.2%
...
Total	24,900,292,747	842,849	562,022	29,543	70.7%

Panel D: Most Valuable Programming Languages

	Total ξ	Mean ξ	Median ξ	N Repos	% Permissive
Python	6,853,266,877	1,154,526	826,076	5,936	75.1%
TypeScript	1,924,509,635	834,928	623,765	2,305	80.2%
JavaScript	1,738,402,866	552,927	272,963	3,144	72.2%
Jupyter Notebook	1,565,510,149	1,302,421	941,121	1,202	79.9%
C#	1,256,454,230	794,721	565,036	1,581	71.8%
Java	1,216,853,117	640,112	414,790	1,901	74.8%
C++	974,203,192	987,035	669,700	987	66.6%
Shell	782,874,993	794,797	555,364	985	73.4%
Go	754,409,476	517,428	228,037	1,458	80.0%
HTML	605,445,891	663,866	363,772	912	59.4%

Table 4
Repository Value by Topic

This table reports summary statistics of repository values by repository topic. ξ is the estimated repository value reported in 2023 dollars. Topic Score, measured between zero and one, reflects how related a repository is to the topic. % Permissive is the percent of repositories with permissive licenses.

	Mean ξ	Median ξ	Total ξ	Mean Topic Score	N Repos	% Permissive
Core AI and ML	880,828	623,525	2,850,360,153	0.60	3,236	70.3%
AI Applications	824,267	584,499	2,853,612,365	0.58	3,462	70.4%
Digital Media	618,571	359,345	1,214,874,229	0.60	1,964	64.2%
Education and Learning	548,422	369,979	1,188,979,071	0.51	2,168	52.1%
Advanced Data Analysis	488,397	330,059	1,933,075,386	0.44	3,958	74.3%
Cloud Infrastructure and DevOps	482,957	360,650	5,319,776,510	0.56	11,015	82.9%
Security	425,193	261,813	1,436,727,161	0.53	3,379	77.2%
Configuration and Templates	370,301	258,140	856,135,980	0.45	2,312	81.8%
Development Tools	361,109	203,698	1,926,875,938	0.48	5,336	75.1%
OS and Platforms	356,741	180,554	682,089,698	0.47	1,912	58.3%
General Data Handling	325,652	211,294	2,492,538,843	0.38	7,654	76.2%
Software Engineering	322,105	206,260	5,649,729,492	0.40	17,540	73.7%
Back-End Web Development	308,481	171,495	587,655,412	0.46	1,905	66.6%
Front-End Web Development	303,391	141,752	700,226,773	0.48	2,308	64.5%
Documentation	277,836	156,394	1,188,581,911	0.39	4,278	67.9%
Community and Governance	275,448	126,018	91,173,314	0.33	331	65.9%

Table 5
Repository Value and Future Popularity

This table reports regression results to validate the measure of private value for repositories (i.e., the dependent variable, ξ). The key variable of interest is $\ln(\text{Stars} + 1)$. “Stars,” as a measure of popularity, is the number of stars as of February 2024. See Table A1 for the definition of variables. Standard errors are double clustered by industry and year in Columns (1) through (4), clustered by firm in Column (5), and are reported in parentheses. ***, **, and * indicate significance at the 1, 5, and 10% levels, respectively.

	(1)	(2)	(3)	(4)	(5)
$\ln(\text{Stars} + 1)$	0.123*** (0.023)	0.107*** (0.022)	0.089*** (0.015)	0.074*** (0.021)	0.066*** (0.015)
$\ln(\text{Mkt Cap})$	1.815*** (0.115)	1.803*** (0.117)	1.734*** (0.064)	1.827*** (0.070)	1.667*** (0.132)
$\ln(\text{Volatility})$	0.406*** (0.056)	0.596*** (0.049)	0.620*** (0.035)	0.463*** (0.022)	
$\ln(\text{Employees})$	-0.289* (0.151)	0.102 (0.146)	0.196*** (0.053)	-0.157*** (0.046)	
$\ln(\text{Total Patent Value} + 1)$	0.113 (0.075)	0.050 (0.048)	0.078** (0.029)	0.211*** (0.028)	
Observations	28,388	28,388	28,388	28,388	28,388
Adj. R ²	0.782	0.800	0.813	0.850	0.858
Year FE	✓	✓			
Industry FE		✓			
Industry x Year FE			✓	✓	
Firm FE				✓	
Firm x Year FE					✓

Table 6
Determinants of Repository Value

This table reports which repository characteristics (Panel A), firm characteristics (Panel B), and product market characteristics (Panel C) are correlated with GitHub repository value (ξ). Control variables include market capitalization, volatility, employees and patent value. See Table A1 for the definition of variables. Standard errors double clustered by industry and year are reported in parentheses. ***, **, and * indicate significance at the 1, 5, and 10% levels, respectively.

<i>Panel A: Repository Characteristics</i>								
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
ln(Stars + 1)	0.090*** (0.017)	0.089*** (0.022)	0.088*** (0.019)	0.054** (0.021)	0.099*** (0.022)	0.085*** (0.021)	0.148*** (0.030)	0.097** (0.033)
Copyleft License	0.089*** (0.019)							0.105** (0.038)
Other License	0.037 (0.027)							0.019 (0.025)
Template		-0.248** (0.096)						-0.178*** (0.046)
Complementarity			-0.302*** (0.071)					-0.221*** (0.048)
Novelty				0.499*** (0.093)				0.379*** (0.092)
ln(Repo Size + 1)					-0.026** (0.008)			-0.014 (0.014)
ln(N Repos + 1)						-0.273** (0.097)		-0.264*** (0.041)
ln(N Issues Opened)							-0.080*** (0.023)	-0.045* (0.022)
Observations	28,690	28,690	28,690	28,690	28,690	28,690	28,690	28,690
Adj. R ²	0.819	0.819	0.821	0.820	0.819	0.822	0.820	0.825
Controls	✓	✓	✓	✓	✓	✓	✓	✓
Industry x Year FE	✓	✓	✓	✓	✓	✓	✓	✓
<i>Panel B: Firm Characteristics</i>								
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
ln(Stars + 1)	0.089*** (0.018)	0.082*** (0.020)	0.085*** (0.022)	0.089*** (0.017)	0.087*** (0.018)	0.086*** (0.021)	0.085** (0.031)	0.079** (0.032)
Market-to-Book	-0.000 (0.020)							-0.008 (0.023)
Return-on-Assets		0.081** (0.027)						0.083** (0.035)
Investment			0.043 (0.056)					-0.001 (0.041)
Return (t-1)				0.002 (0.015)				0.005 (0.017)
Sales Growth					0.043 (0.036)			0.032 (0.034)
Tangibility						0.039 (0.068)		-0.028 (0.050)
R&D Exp/Total Assets							0.089 (0.070)	0.090 (0.060)
R&D Exp Missing							0.518*** (0.136)	0.472*** (0.096)
Observations	27,583	27,583	27,583	27,583	27,583	27,583	27,583	27,583
Adj. R ²	0.806	0.807	0.806	0.806	0.806	0.806	0.807	0.808
Controls	✓	✓	✓	✓	✓	✓	✓	✓
Industry x Year FE	✓	✓	✓	✓	✓	✓	✓	✓

Panel C: Product Market Characteristics

	(1)	(2)	(3)	(4)	(5)	(6)
ln(Stars + 1)	0.068*** (0.013)	0.078*** (0.014)	0.065*** (0.013)	0.070*** (0.013)	0.069*** (0.014)	0.064*** (0.014)
Market Power	0.102*** (0.018)					0.067*** (0.019)
Product Market Centrality		0.017 (0.011)				0.135*** (0.022)
Scope			-0.149*** (0.018)			-0.100*** (0.019)
Product Market Similarity				-0.088 (0.047)		-0.010 (0.032)
Product Market Fluidity					-0.158*** (0.022)	-0.131** (0.044)
Observations	23,735	23,735	23,735	23,735	23,735	23,735
Adj. R ²	0.824	0.821	0.825	0.822	0.824	0.827
Controls	✓	✓	✓	✓	✓	✓
Industry x Year FE	✓	✓	✓	✓	✓	✓

Table 7
Repository Values and Firm Output

This table reports the relation between the value of all repositories posted by a firm in year t (ξ) and the firm's future outcomes over horizons from year $t + 1$ to year $t + 3$, as formulated in Equation 2. Dependent variables include the growth of: sales, profits, number of employees, number of patents, and patent value. Control variables include one lag of the dependent variable, market capitalization, idiosyncratic volatility, employees, and patent value. See Table A1 for the definition of variables. Standard errors double clustered by industry and year are in parentheses. ***, **, and * indicate significance at the 1, 5, and 10% levels, respectively.

	Firm (Horizon)		
	1	2	3
<i>Panel A: Sales</i>			
$\ln(\xi + 1)$	0.002* (0.001)	0.005*** (0.001)	0.008*** (0.002)
<i>Panel B: Profits</i>			
$\ln(\xi + 1)$	0.001* (0.001)	0.004** (0.001)	0.007** (0.002)
<i>Panel C: Labor</i>			
$\ln(\xi + 1)$	0.001 (0.001)	0.004** (0.001)	0.008*** (0.002)
<i>Panel D: Value of Patents</i>			
$\ln(\xi + 1)$	0.005*** (0.001)	0.013** (0.004)	0.021*** (0.004)
<i>Panel E: Number of Patents</i>			
$\ln(\xi + 1)$	0.005*** (0.001)	0.011*** (0.002)	0.016*** (0.003)
Controls	✓	✓	✓
Industry FE	✓	✓	✓
Year FE	✓	✓	✓

Appendices

Appendix A

Table A1
Variable Definitions

Variable	Definition
Complementarity	Score between zero and one that measures how much the repository complements the firm's commercial products (ChatGPT).
Copyleft License	Indicator variable that equals one if the repository has a license with some copyleft restrictions (GitHub API).
Employees	Number of employees (Compustat).
GitHub	Indicator variable that equals one after the firm releases its first repository (GHArchive).
Investment	CAPX scaled by lagged total assets (Compustat).
Market Capitalization	Share price times the number of shares outstanding (CRSP).
Market Power	An estimate of markups assuming constant returns to scale developed by (Pellegrino, 2024).
Market-to-book	Ratio of market capitalization to book equity, where book equity is calculated following Davis et al. (2000) (CRSP, Compustat).
N Commits	Number of commits across all repositories owned by the firm in that month (GHArchive).
N Issues Opened	Cumulative number of issues opened for a repository as of December 31, 2023 (GHArchive).
N Repos (t)	Cumulative number of repositories released by a firm prior to month t (GHArchive).
Novelty	Score between zero and one that measures how novel or groundbreaking a repository is compared to existing solutions, focusing on whether it introduces new ideas, techniques, or approaches (ChatGPT).
Number of Patents	Number of patents granted (Kogan et al., 2017).
Other License	Indicator variable that equals one if the repository has a customized license that cannot be cleanly classified into "copyleft" or "permissive" categories.
Patent Value	An estimate of the economic value of patents using stock market returns around the patent grant date (Kogan et al., 2017)
Product Market Centrality	Eigenvector centrality calculated from a network created by product market similarity scores (Hoberg and Phillips, 2016).
Product Market Fluidity	A measure of how intensively the product market around a firm is changes (Hoberg et al., 2014).

Continued on the next page

Continued

Variable	Definition
Product Market Similarity	A measure of how similar a firm's products are to its peers', from Hoberg and Phillips (2016) (Hoberg-Phillips Data Library).
Profits	Sale minus COGS, deflated by the CPI (Compustat)
R&D Exp/Total Assets	Research and development expense scaled by lagged total assets (Compustat).
R&D Exp Missing	Indicator variable equal to one if R&D expense is missing (Compustat).
Repo Size	Byte size of a repository as of February 2024 (GitHub API).
Return (t)	Returns from month t (CRSP).
Return-on-Assets	Net income divided by lagged total assets (Compustat).
Sales	Annual Sales (Compustat).
Scope	Number of industries in which the firm operates, see Hoberg and Phillips (2024) (Hoberg-Phillips Data Library).
Stars	Number of stars of a repository as of February 2024 (GitHub API).
Tangibility	Property, plant, and equipment scaled by total assets (Compustat).
Template	Indicator variable that equals one if the repository is configured as a template, which allows copies to be created without retaining the commit history (GitHub API).
Topic Score	Score between zero and one that measures how much the repository relates to the given topic (ChatGPT).
(Idiosyncratic) Volatility	Standard deviation of daily returns over one month. Idiosyncratic volatility is similarly defined using returns net of market returns (CRSP).
ξ	An estimate of the economic value of repositories (in 2023 dollars) using stock market returns around the repository release date.

Internet Appendix A

Open-Source Licenses

Table IA1
License Classification Based on Permission Levels

This table classifies common licenses of GitHub repositories based on their permission levels. We adopt the Open Source Initiative’s definition of open source, which stipulates that an open source license must not discriminate against any person, restrict other software, or be specific to a product, among other criteria. See <https://opensource.org/osd/> for details. Among open source licenses, permissive licenses impose minimal restrictions, whereas copyleft licenses require that derived works also be open source.

Type of licenses	Restrictions	Benefits	Costs	Examples
Permissive License				
Open source and permissive	Keep copyright information	Compatibility with both open source and proprietary projects	Limited protection of the original developer’s work	MIT, Apache 2.0
Copyleft License				
Open source and weak copyleft	Copyleft for the original codes	Encourages contributions to the open source component while permitting proprietary integration	Incompatible with proprietary projects with static linking	LGPL
Open source and strong copyleft	Can use but derivative work must also be open source	Preserve the open nature; Monetize	Incompatible with proprietary projects and compliance burden	GPL, AGPL
Other License				
Source available but limit use in certain products	Limits the use of the software in specific commercial or competitive scenarios.	Monetize software by offering additional commercial licenses/through complementary products	Less contribution from both individual users and commercial users	Amazon Software License
Source available but limit commercial use	Only for non-commercial purposes	Encourages contributions and community engagement while protecting against certain commercial uses	Close monitoring and compliance enforcement; Less contribution from commercial users	CC-BY-NC-4.0
No license	By default illegal to use, distribute, or modify the code	IP protection, flexibility in choosing license later	Discourage adoption and contribution	
Open source with copyright notice	Can be either permissive or copyleft			

Internet Appendix B

This appendix provides details on the steps we take to classify or evaluate GitHub repositories based on their topics, the extent to which they complement firms' commercial offerings, and novelty using large language models (LLMs).

Classifying GitHub Repositories by Topic

While repository admins can add topics to increase the visibility of their projects, there are no specific requirements regarding which topics they can use. These topic labels can be assigned based on the intended purpose, subject area, affinity groups, or other important qualities. Therefore, the potential choice of topics is unlimited. In addition, many repositories remain unlabeled. To address these challenges, we rely on the GitRanking taxonomy, proposed by [Sas et al. \(2023\)](#), to limit the topic space for repository classification based on a structured set of topics. Additionally, we employ LLMs to infer appropriate topics for repositories by analyzing the content of the repositories and identifying relevant themes, even when explicit topics are not provided.

GitRanking is a taxonomy consisting of 301 labels derived from 121,000 GitHub topics. These 301 topic labels are organized into distinct levels based on their meanings. To balance the breadth and specificity of topics for our purposes, we elect to use the third level, which is comprised of 62 topics. To further refine these topics, we use ChatGPT to group similar topics, enhance clarity, and reduce overlap between topics. This process resulted in 17 distinct topics, which we use to classify the repositories in our sample. Definitions of these topics are provided in the prompt below.

Next, we use ChatGPT to assess how closely a repository aligns with a given topic, assigning a relatedness score ranging from 0 to 1. The information provided to ChatGPT includes the repository's name, description, main programming language, self-reported topics, and website, all obtained via the GitHub API. We explicitly allow ChatGPT to incorporate

external knowledge beyond the provided data. Furthermore, to ensure consistency in evaluating a repository’s relatedness to a topic, we developed a scoring reference and repeatedly tested it on a small subsample of repositories, confirming that score variations typically stay within 0.1.

Specifically, we use the following model parameters and prompt for OpenAI’s API:

Model: gpt-4o-2024-08-06

Temperature: 0

Seed: 2024

Prompt: Assign relevance weights (0 to 1) to predefined topic labels for a GitHub repository. The weight represents how relevant each topic is to the repository. Your weights should follow the scoring reference and definitions of topic labels as below.

Scoring Reference:

- **0.0:** The topic is **completely irrelevant** to the repository. There is **no code, documentation, or features** associated with this topic. The topic does not apply in any way.
- **0.1:** The topic has **extremely low relevance**. It is **mentioned only once** or in a **minor part** of the repository (e.g., a reference in a single file or a brief mention in documentation). There is **no meaningful functionality** related to this topic.
- **0.2:** The topic has **very low relevance**. There is a **small, secondary feature** related to the topic, but it plays a **minimal role** in the repository. It is **not central** to the repository’s purpose and may only be used in a **supporting or optional capacity**.
- **0.3:** The topic has **low relevance**. The repository includes **some functionality** or content related to the topic, but it is a **minor or auxiliary component**. The topic is **not integral** to the primary goals of the repository, though it may be referenced or used in **specific parts** of the project.

- **0.4:** The topic has **below-average relevance**. It plays a **noticeable role** in the repository, but it is **not essential**. There are **clear references, code, or features** related to the topic, but they are **not a major focus**.
- **0.5:** The topic has **moderate relevance**. It is one of **several key areas** covered by the repository. A **substantial portion** of the repository's code, features, or documentation relates to this topic, but it is **not the main focus** of the project.
- **0.6:** The topic has **moderately high relevance**. The topic is one of the **core areas** of the repository, with a **significant portion** of the code, features, or documentation focused on it. The repository relies heavily on this topic, but other topics are also important.
- **0.7:** The topic has **high relevance**. A **large portion** of the repository is built around this topic. The topic plays a **central role** in the repository's features, functionality, or design. Most of the repository's content is directly related to this topic, but there are still other areas of focus.
- **0.8:** The topic has **very high relevance**. It is a **primary focus** of the repository, with **most features, code, and documentation** centered around it. Nearly all content is related to this topic, with only a few secondary areas.
- **0.9:** The topic has **near-perfect relevance**. The repository is **almost entirely centered** around this topic. The vast majority of its code, design, and purpose are related to this field, with **only minor mentions** of other topics.
- **1.0:** The topic has **perfect relevance**. The repository's **sole purpose** is to serve this topic. **Every feature, line of code, and piece of documentation** is directly related to this topic, with **no other topics** playing a significant role.

Topic labels to assign weights to:

- **Digital Media:** Projects related to game development, video games, animation, camera software, image processing, audio or video editing, and interactive media.
- **General Data Handling:** Projects focused on managing data, including data structures, file systems, databases, and ETL (Extract, Transform, Load) processes.
- **Advanced Data Analysis:** Projects involving complex data analysis, such as big data, bioinformatics, data science, text analysis, time series analysis, or data visualization.
- **Security:** Projects related to cryptography, data protection, authentication, network security, privacy, or detecting threats like phishing.
- **Cloud Infrastructure and DevOps:** Projects centered on cloud computing, CI/CD pipelines, distributed computing, microservices, backups, infrastructure automation, or serverless computing.
- **Software Engineering:** Projects focused on software architecture, testing, program analysis, design patterns, or software development methodologies like Agile or Scrum.
- **Front-End Web Development:** Projects involving client-side development, including UI/UX design, HTML/CSS, and JavaScript frameworks.
- **Back-End Web Development:** Projects centered on server-side development, including APIs, databases, and routing.
- **Core AI/ML:** Projects focused on foundational machine learning concepts such as neural networks, deep learning, semi-supervised learning, or reinforcement learning.
- **AI Applications:** Projects applying AI techniques in specific domains, such as robotics, computational biology, computer vision, or natural language processing (NLP).

- **Development Tools:** Projects related to building or maintaining programming tools like compilers, interpreters, validators, debugging tools, IDEs, or version control systems.
- **Operating Systems and Platforms:** Projects focused on OS development, kernel development, embedded systems, or platform-specific development (e.g., Windows, Linux, Android).
- **Documentation:** Projects primarily providing manuals, guides, API documentation, or technical standards.
- **Community and Governance:** Projects focused on open-source community guidelines, such as codes of conduct, contribution guidelines, or governance policies.
- **Education and Learning:** Projects designed for educational purposes, such as tutorials, training modules, or coding bootcamps.
- **Configuration and Templates:** Projects offering pre-configured setups, boilerplate code, or templates for quick project initialization (e.g., Dockerfiles, CI/CD configs).
- **Other:** For any project that doesn't fit the above categories. The weight should be zero if no miscellaneous topics are relevant to the repository, or a non-zero value if there are other significant topics.

Use both the repository details and relevant knowledge you have about this repository to make your decision.

Evaluating Complementarity of GitHub Repositories

Similarly, we use LLMs to evaluate the complementarity of GitHub repositories to their owners' commercial products, which is defined as directly supporting or enhancing the firm's core business products. In addition to the repository's name, description, main programming

language, self-reported topics, and website, we also include the Compustat name of the firm owning the repository. We explicitly allow ChatGPT to incorporate external knowledge beyond the provided data, and include a scoring reference. We have ChatGPT not only provide the complementarity score but also the name of the commercial product that this repository complements, with which we confirm the validity of the scores.

To illustrate the intuition of the resulting complementarity scores, we provide the following examples. The first example is the repository “WhatsApp/StringPacks,” which is a library designed to store translation strings in a more efficient binary format for Android applications. This library can operate completely independently of WhatsApp’s messaging services and be used by any relevant Android application. This repository has a complementarity score of 0. In contrast, the second example is the repository “WhatsApp/stickers,” which contains iOS and Android sample apps as well as an API for creating third-party sticker packs for WhatsApp. This project directly complements the WhatsApp messaging app, resulting in a complementarity score of 0.8.

We use the following model parameters and prompt for OpenAI’s API:

Model: gpt-4o-2024-08-06

Temperature: 0

Seed: 2024

Prompt:

Evaluate whether a GitHub repository owned by a US public firm complements the firm’s commercial products or operates as a standalone project based on the following scoring reference. Use both the repository details and relevant knowledge you have about this repository to make your decision. Only return in JSON and do not include anything else. In your JSON response, include the following fields: `repo_id`, `comp_score` (complementarity score), and `comm_product` (the commercial product to which the repository complements).

Scoring Reference:

- **0.0:** The repository is entirely standalone. It has no overlap or connection with any

of the firm's commercial products. It operates independently of the company's core business.

- **0.1:** The repository shows minimal potential relevance or use in conjunction with the firm's commercial products but is not designed for or marketed as part of the firm's offerings.
- **0.2:** The repository may have slight overlaps with the firm's commercial offerings but is not positioned as a key integration. It could potentially be used with the firm's products but has no direct integration or clear marketing as a complementary tool.
- **0.3:** The repository shows some potential to complement the firm's products but is still largely standalone. There might be some integrations, but they are not essential or exclusive to the firm's ecosystem.
- **0.4:** The repository provides a small but noticeable enhancement to the firm's commercial products. However, the connection to the commercial offering is weak, and the repository is still usable independently.
- **0.5:** The repository offers some clear value to the firm's commercial products but is not a core or exclusive component. It may integrate with or enhance the firm's product, but its relevance is moderate.
- **0.6:** The repository offers strong support for the firm's products and likely exists to enhance or complement the product experience. However, it is still not fully dependent on the commercial product.
- **0.7:** The repository is closely linked with the firm's commercial product and provides substantial enhancements or integrations. It is marketed or documented as a useful component for customers of the firm's product.

- **0.8:** The repository strongly complements the firm’s commercial product and is used almost exclusively within the context of that product. However, it may still be used in other contexts with significant effort.
- **0.9:** The repository is nearly indispensable for customers using the firm’s commercial product. It is closely tied to the product, and its functionality is largely dependent on it.
- **1.0:** The repository is entirely and exclusively built to complement and support the firm’s commercial product. It has no utility outside of the firm’s product and is crucial for its full use. The repository cannot function independently and exists solely to enhance the firm’s commercial offering.

Category-Specific Definitions:

- **Complementary Repositories:** These repositories directly support or enhance the firm’s core business offerings.
- **Standalone Repositories:** These are open-source projects, tools, or experiments that do not contribute to or enhance the firm’s main commercial products, even if created or maintained by the firm.

Evaluating Novelty of GitHub Repositories

Lastly, we use LLMs to evaluate the novelty of GitHub Repositories. To do so, we provide the repository’s name, description, main programming language, self-reported topics, and website. We explicitly allow ChatGPT to incorporate external knowledge beyond the provided data, and include a scoring reference.

Again, to illustrate the intuition for the resulting novelty scores, we consider the following examples. First, the average novelty score for repositories affiliated with the “LinkedInLearning” organization account is 0.1. These projects are often exercises associated with courses on

the LinkedIn Learning platform. In contrast, the repository “google-deepmind/alphafold”, which hosts the open-source code of AlphaFold (an AI system developed by the 2024 Chemistry Nobel Prize winners that predicts a protein’s 3D structure), has a novelty score of 0.8, the highest in our sample.^{IA1}

We use the following model parameters and prompt for OpenAI’s API:

Model: gpt-4o-2024-08-06

Temperature: 0

Seed: 2024

Prompt:

Evaluate the originality of a GitHub repository. Originality measures how novel or groundbreaking a repository is compared to existing solutions, focusing on whether it introduces new ideas, techniques, or approaches.

Use both the repository details and relevant knowledge you have about this repository to make your decision. Use the scoring reference provided for consistency. Only return in JSON and do not include anything else.

Originality Scoring Reference:

- **0.0:** The repository introduces no new ideas or techniques. It is a near-complete replication of existing solutions with no modifications.
- **0.1:** The repository shows minimal originality, with minor tweaks or adaptations of well-established methods, but still follows existing patterns closely.
- **0.2:** Low originality. The repository includes slight variations or small improvements on existing solutions but doesn’t introduce any novel concepts.
- **0.3:** The repository adds some new ideas or features, but these are incremental and build directly on existing work.

^{IA1} Note that the Nobel Prize announcement occurred after the model’s training period, and therefore the score was not influenced by the Nobel news.

- **0.4:** Below-average originality. The repository offers a combination of existing techniques with minor innovations or optimizations.
- **0.5:** Moderate originality. The repository introduces some interesting and unique features or techniques, but they are not highly groundbreaking.
- **0.6:** The repository shows notable originality. It provides a fresh approach to solving a problem, though the idea may not be entirely new.
- **0.7:** High originality. The repository demonstrates a new concept or method that has the potential to influence other projects or domains.
- **0.8:** Very high originality. The repository presents a significantly novel approach, introducing new methodologies, techniques, or tools that are not widely available.
- **0.9:** Nearly groundbreaking. The repository offers a unique solution that stands out as highly innovative compared to others in the field.
- **1.0:** Completely original. The repository introduces entirely new concepts, techniques, or tools that could redefine the domain or set new standards.

Internet Appendix C

Estimating Repository Value

This appendix provides an extended description of our procedure to estimate repository value. Further details and discussion of assumptions can be found in [Kogan et al. \(2017\)](#).

The procedure involves observing stock returns in the three-day window following the announcement of the repository, $[t, t+2]$. We choose this window for multiple reasons. First, it is the same window used by prior studies, and so ensures the comparability of our estimates with those for other assets. Second, it limits the probability of other events contaminating the estimates. While expanding the window could capture more of the market reaction to minor repositories that do not induce a significant or immediate reaction from investors upon announcement, the additional noise could render the estimates largely uninformative. Third, we contend that announcements of open-source repositories are predominantly unexpected, so we do not include days prior to the announcement in the announcement window.^{IA2}

To begin our estimation procedure, we remove fluctuations in daily returns attributable to market movements by subtracting the market return from each firm’s daily return. We then cumulate these market-adjusted returns over the three-day announcement window for repository i , which we label R_i . We assume that R_i is a function of both investor reaction to the repository announcement, v_i , and idiosyncratic noise, ε_i , such that

$$R_i = v_i + \varepsilon_i. \tag{3}$$

We construct the estimate of repository value as the product of the investor reaction to

^{IA2} In estimating patent values, [Kogan et al. \(2017\)](#) apply an adjustment for the fact that the announcements are of patent *grants*, while information regarding the patent is first revealed to investors when the patent application is filed. The market reaction on the grant date therefore reflects only a portion of the value corresponding to the resolution of uncertainty as to whether the patent is granted. In the context of GitHub repositories, however, this adjustment is not needed. Information about repositories is not systematically shared prior to the repositories being open-sourced, and so market reactions within the announcement windows reflect the full value of the repositories.

the repository announcement and the firm’s market capitalization on the day prior to the announcement. If multiple repositories are announced on the same day, we assume the value is evenly distributed across those repositories. Given that repository announcements do not follow a typical schedule,^{IA3} multiple repository announcements on the same day tend to, anecdotally, correspond to a single project. The value of repository i , ξ_i , is thus calculated as

$$\xi_i = \frac{1}{N_i} E[v_i | R_i] M_i, \quad (4)$$

where N_i is the number of repositories announced on that day, $E[v_i | R_i]$ is the expected return attributable to the repository announcement conditional on observing the three-day cumulative market-adjusted return R_i , and M_i is the market capitalization of the firm on the day prior to the repository announcement.

To estimate the conditional expected return in Equation (4), we adopt the same distributional assumptions about v and ε as Kogan et al. (2017).^{IA4} Note that the distributional assumption regarding v_i implies that repositories have strictly positive values. While it is possible that open-source projects provide value to competitors that make the projects less valuable to the firm itself, we assume that firms will only choose to make projects open source if the net effect still results in a positive value for the firm. Under these assumptions, the conditional expected return can be calculated as

$$E[v_i | R_i] = \delta R_i + \sqrt{\delta} \sigma_{\varepsilon ft} \frac{\phi\left(-\sqrt{\delta} \frac{R_i}{\sigma_{\varepsilon ft}}\right)}{1 - \Phi\left(-\sqrt{\delta} \frac{R_i}{\sigma_{\varepsilon ft}}\right)}, \quad (5)$$

where ϕ and Φ represent the standard normal PDF and CDF, respectively, and δ denotes

^{IA3} In comparison, patent grants are announced every Tuesday.

^{IA4} Specifically, we assume v_i follows a normal distribution truncated at zero such that $v_i \sim \mathcal{N}^+(0, \sigma_{v ft}^2)$ and ε_i follows a normal distribution such that $\varepsilon_i \sim \mathcal{N}(0, \sigma_{\varepsilon ft}^2)$. Thus, both distributions vary across firms, f , and time, t .

the signal-to-noise ratio,

$$\delta = \frac{\sigma_{vft}^2}{\sigma_{vft}^2 + \sigma_{\varepsilon ft}^2}. \quad (6)$$

We adopt the same simplifying assumption as [Kogan et al. \(2017\)](#) that δ is the same for all firms and all time periods. We believe this assumption is reasonable in our setting due to the relatively short time period, which begins in 2015. This assumption still allows σ_{vft}^2 and $\sigma_{\varepsilon ft}^2$ to vary across firms and time, but only in constant proportion. To estimate δ , we compare the variance of returns in the announcement window to that of returns over other three-day periods for the same firm within the same year. This comparison takes the regression form

$$\ln(R_{fd}^2) = \gamma I_{fd} + \lambda_{dow} + \eta_{fy} + u_{fd}, \quad (7)$$

where R_{fd} is the three-day cumulative market-adjusted return for firm f on day d , I_{fd} is an indicator variable that equals one if there is a repository announcement by firm f on day d , λ_{dow} are day-of-week fixed effects, and η_{fy} are firm-year fixed effects. Importantly, this regression only includes firms that have a repository announcement at some point in the sample period. The estimated $\hat{\delta}$ can be calculated from the resulting estimate $\hat{\lambda}$ as $\hat{\delta} = 1 - e^{-\hat{\gamma}}$. For our main sample of repositories with available public dates, $\hat{\gamma} = 0.0359$ and $\hat{\delta} = 0.0353$.

Finally, we estimate $\sigma_{\varepsilon ft}^2$ for each firm within each year as

$$\sigma_{\varepsilon ft}^2 = \frac{3\sigma_{ft}^2}{1 + 3d_{ft}(e^{-\hat{\gamma}} - 1)}, \quad (8)$$

where d_{ft} is the fraction of days in the given year that are announcement days for firm f and σ_{ft}^2 is the variance of daily market-adjusted returns calculated within each firm for each year.

Internet Appendix D

Correlations

This appendix reports univariate correlations among all pairs of variables included in our analysis of the determinants of repository value. The correlation matrix is presented in Table [IA2](#).

Table IA2
Correlation Matrix

This table presents a correlation matrix of all variables included in our analysis of the determinants of repository value. Each variable is defined, along with its data source, in Appendix A1.

	ln(ξ)	ln(Stars + 1)	ln(Mkt Cap)	ln(Volatility)	ln(Employees)	ln(Patent Port ξ + 1)	Permissive License	Other License	Template	Complementarity
ln(Stars + 1)	0.251									
ln(Mkt Cap)	0.851	0.215								
ln(Volatility)	-0.263	-0.181	-0.507							
ln(Employees)	0.689	0.081	0.855	-0.422						
ln(Patent Port ξ + 1)	0.688	0.209	0.839	-0.587	0.768					
Permissive License	0.156	0.079	0.214	-0.072	0.233	0.136				
Other License	-0.143	-0.072	-0.196	0.066	-0.224	-0.124	-0.966			
Template	-0.010	-0.007	0.003	0.010	-0.018	-0.018	-0.004	0.007		
Complementarity	-0.020	-0.013	0.069	-0.019	0.109	0.062	0.152	-0.147	0.020	
Novelty	0.189	0.432	0.111	0.052	0.060	0.094	0.090	-0.087	-0.039	-0.130
ln(Repo Size + 1)	0.021	0.345	0.021	-0.049	-0.022	0.054	-0.032	0.025	-0.019	0.123
ln(N Repos + 1)	0.556	0.035	0.699	-0.327	0.583	0.549	0.181	-0.170	0.039	0.086
ln(Issues Opened + 1)	0.064	0.727	0.060	-0.160	-0.026	0.092	0.075	-0.071	-0.018	0.167
Market-to-Book	0.072	-0.103	0.065	0.212	0.007	-0.112	0.049	-0.045	0.007	0.116
Return-on-Assets	0.503	0.189	0.551	-0.380	0.409	0.615	0.024	-0.018	-0.029	-0.099
Investment	0.402	0.050	0.487	0.051	0.641	0.330	0.216	-0.201	-0.016	0.095
Return (t-12 to t-1)	0.189	0.026	0.168	0.039	-0.016	-0.035	0.020	-0.017	0.016	0.025
Sales Growth	0.096	0.053	0.018	0.360	-0.047	-0.252	0.085	-0.071	0.006	0.034
Tangibility	0.424	-0.009	0.507	0.052	0.718	0.367	0.217	-0.203	-0.009	0.077
R&D Exp/Total Assets	0.083	0.014	0.033	0.417	0.105	-0.100	0.123	-0.112	-0.020	0.209
R&D Exp Missing	-0.193	-0.081	-0.218	-0.002	-0.101	-0.220	-0.036	0.023	-0.014	-0.170
Market Power	-0.154	0.092	-0.244	0.113	-0.449	-0.189	-0.124	0.126	0.007	-0.143
Scope	-0.059	0.034	-0.040	-0.221	-0.312	0.000	-0.054	0.055	0.049	-0.016
ln(PM Centrality)	-0.178	-0.020	-0.216	0.002	-0.190	-0.178	-0.029	0.017	-0.006	-0.008
PM Similarity	-0.192	-0.020	-0.191	-0.099	-0.363	-0.141	-0.073	0.065	0.020	-0.026
PM Fluidity	-0.103	0.157	-0.093	-0.289	-0.190	-0.008	0.004	-0.007	-0.002	0.032

	Novelty	ln(Repo Size + 1)	ln(N Repos + 1)	ln(Issues Opened + 1)	Market-to-Book	Return-on-Assets	Investment	Return (t-12 to t-1)
ln(Repo Size + 1)	0.148							
ln(N Repos + 1)	0.047	0.009						
ln(Issues Opened + 1)	0.263	0.378	-0.033					
Market-to-Book	-0.040	-0.023	0.096	-0.110				
Return-on-Assets	0.151	0.055	0.299	0.079	-0.097			
Investment	0.137	-0.064	0.329	-0.065	0.129	0.172		
Return (t-12 to t-1)	0.028	0.000	0.118	-0.014	0.268	0.018	0.037	
Sales Growth	0.099	-0.045	0.016	-0.027	0.326	-0.131	0.384	0.246
Tangibility	0.117	-0.078	0.350	-0.120	0.147	0.170	0.892	0.017
R&D Exp/Total Assets	0.087	-0.028	0.016	-0.052	0.382	-0.215	0.527	0.158
R&D Exp Missing	-0.029	-0.041	-0.334	-0.052	-0.154	-0.060	-0.150	-0.102
Market Power	0.092	0.057	-0.111	0.106	0.024	0.098	-0.214	-0.087
Scope	-0.090	0.080	0.071	0.096	-0.056	0.053	-0.555	0.128
ln(PM Centrality)	-0.066	-0.005	-0.267	0.031	-0.107	-0.204	-0.328	-0.089
PM Similarity	-0.092	0.058	-0.045	0.081	-0.065	-0.215	-0.495	-0.018
PM Fluidity	-0.077	0.047	-0.151	0.197	-0.191	-0.162	-0.442	-0.022

	Sales Growth	Tangibility	R&D Exp/Total Assets	R&D Exp Missing	Market Power	Scope	ln(PM Centrality)	PM Similarity
Tangibility	0.258							
R&D Exp/Total Assets	0.520	0.487						
R&D Exp Missing	-0.125	-0.075	-0.342					
Market Power	0.130	-0.332	-0.073	-0.111				
Scope	-0.227	-0.571	-0.327	-0.017	0.145			
ln(PM Centrality)	-0.180	-0.241	-0.064	0.261	-0.031	0.271		
PM Similarity	-0.183	-0.558	-0.305	0.015	0.212	0.594	0.375	
PM Fluidity	-0.267	-0.431	-0.177	0.112	-0.003	0.555	0.668	0.461