

# Forks at Midnight: The Diffusion of Open-Source AI Across Developers and Firms

*Anders Holm\**, *Femi Adebayo*

Generative Economic Research Institute (GERI) • Center for AI and Knowledge Work (CAIKW)

Generative Economic Review • May 18, 2026

GER 1.11

---

**JEL Classification:** O33, L86, M15, O31, C81

**Keywords:** open source software, artificial intelligence, GitHub, technology diffusion, large language models, Bass diffusion model, structural break, developer ecosystems, langchain, generative AI, organizational learning, platform-mediated innovation

---

## Abstract

We document the temporal trajectory of open-source artificial intelligence project creation on GitHub from 2018 through 2026 using the GitHub Search API, and we provide the first formal diffusion-model estimation applied to the developer-ecosystem adoption of generative AI. For repositories created in each calendar year we count those tagged with six AI-specific topics (`llm`, `large-language-model`, `generative-ai`, `openai`, `langchain`, `transformer`) and compare against a baseline of all repositories tagged `machine-learning`. Three empirical findings are central. First, repositories tagged with AI-specific topics grew from 260 in 2018 to 48{,}003 in 2026, a compound annual growth rate of 92 percent per year, approximately ten times the corresponding growth rate of the broader machine-learning category (9 percent per year). Second, the inflection in the trajectory is unambiguous and aligns with the public release of large language models: AI-specific repositories grew 7.8-fold between 2022 (2{,}129 repositories) and 2023 (16{,}618 repositories), an order-of-magnitude acceleration concentrated in a single calendar year. A Chow structural-break test rejects the null of parameter stability at the 2022–2023 boundary ( $F = 47.3$ ,  $p < 0.001$ ), and a two-parameter Bass diffusion model fit to the pre-2022 trajectory under-predicts the 2023 count by a factor of 6.1, confirming that the discontinuity lies outside any smooth diffusion envelope calibrated on the prior trajectory. Third, the composition of AI repositories has shifted substantially toward generative-AI-specific topics: the `llm` topic alone grew at a compound annual rate of 131.5 percent over the sample, and `langchain`, a framework released in October

2022, grew from 7 repositories in 2018 to 5{,}155 in 2026. We validate these counts through a tag-accuracy audit of 300 sampled repositories (89 percent true-positive rate) and through overlap-trend analysis showing that the year-over-year overlap rate is approximately stationary, so that double-counting does not materially bias the growth-rate estimates. The paper concludes with three sets of implications: for diffusion theory, where the punctuated-equilibrium pattern we document challenges the smooth S-curve assumption of canonical Bass models; for management practice, where the speed of developer-ecosystem adoption implies that firms delaying AI integration face compounding capability gaps; and for research methodology, where the magnitude of the diffusion poses a challenge for retrospective management research designs whose data-collection timelines lag the phenomenon they study.

---

## 1. Introduction

The diffusion of new technologies through the productive economy is one of the central objects of empirical management research. From the canonical studies of hybrid corn adoption (Griliches, 1957), color television (Bass, 1969), and the early personal computer (Bresnahan and Greenstein, 1996) to the contemporary literature on cloud computing, mobile internet, and artificial intelligence, the methodological pattern is similar: identify an adoption indicator, track its temporal trajectory, decompose into firm-level, sectoral, and geographic components, and fit one of a small number of standard diffusion curves to the trajectory.

### 1.1 The framing hypothesis

This paper makes one central empirical claim. The diffusion of generative artificial intelligence through the open-source developer ecosystem, as measured by GitHub repository creation tagged with AI-specific topics, has occurred at a pace that is faster than the analogous diffusion of any prior technology for which comparable indicators exist. The compound annual growth rate of AI-specific repositories over 2018–2026 is approximately 92 percent per year, approximately ten times the contemporary rate for the broader machine-learning category and approximately twice the historical rate for the diffusion of containerization (Docker repositories, 2014–2018) or microservices architectures (2014–2018). We go beyond this descriptive claim in three ways: we formally estimate the parameters of a Bass diffusion model on the pre-2022 and full-sample trajectories, we conduct a Chow structural-break test at the 2022–2023 boundary, and we validate the tag-based counts through a manual accuracy audit. The empirical pattern has direct implications for the retrospective measurement of technology diffusion in management research, for the design of policy interventions targeting AI literacy and workforce readiness, and for the organizational strategies

of firms navigating the transition.

## 1.2 Four contributions

The paper makes four substantive contributions to the literature on technology diffusion and open-source software.

First, we provide the most comprehensive empirical record of open-source AI project creation across a nine-year window (2018–2026), using publicly accessible GitHub Search API data. The record covers the full range from pre-LLM machine-learning era through the contemporary generative-AI era, with the 2022Q4 capability shock at the center of the window. While the qualitative direction of the trend—rapid growth in AI activity—is widely recognized, the precise magnitude, timing, and composition of the growth have not been formally documented with this specificity. The distinction matters: prior informal estimates of AI repository growth have ranged from ‘50 percent per year’ to ‘exponential,’ a spread too wide to anchor subsequent quantitative work. Our 92 percent CAGR, validated through multiple robustness specifications and a tag-accuracy audit, provides the precise empirical baseline the literature currently lacks.

Second, we document a 7.8-fold acceleration in AI-specific repository creation between 2022 and 2023, a temporal discontinuity that aligns precisely with the public release of large language models. The discontinuity is empirically unambiguous: a Chow structural-break test rejects parameter stability at the boundary ( $F = 47.3$ ,  $p < 0.001$ ), and a Bass diffusion model fit to the pre-2022 trajectory under-predicts the 2023 realization by a factor of 6.1. The formal tests confirm what the raw data suggest: the 2022–2023 transition is not a continuation of the prior trend but a regime change.

Third, we decompose the AI-repository universe into six topic-specific sub-trajectories (llm, large-language-model, generative-ai, openai, langchain, transformer). The langchain topic specifically—corresponding to a framework released in October 2022—illustrates the speed of practitioner-level tool adoption: 7 repositories in 2018 (essentially zero) to 5{,}155 in 2026.

Fourth, we benchmark the AI diffusion trajectory against the comparable trajectories of three prior technology diffusion episodes for which GitHub-based indicators exist: Docker containerization (2014–2018), Kubernetes orchestration (2015–2019), and TensorFlow machine-learning frameworks (2015–2019). The contemporary AI diffusion is approximately twice as fast as any of these comparison cases. We further benchmark the estimated Bass-model innovation parameter ( $p$ ) against the cross-technology distribution reported in Sultan et al. (1990), finding that the AI diffusion’s estimated  $p$  exceeds the 95th percentile of all previously documented diffusion episodes.

### 1.3 Intellectual history of the question

The question of technology diffusion has been addressed empirically through four broad measurement traditions. The first—canonical in agricultural and industrial economics—uses firm-level or farm-level adoption surveys (Griliches, 1957; Comin and Hobijn, 2010). The second—prominent in management research—uses patent filings as adoption proxies (Cockburn et al., 2019). The third—used in the IT and software diffusion literature—uses installed-base indicators (operating-system shares, software-sales receipts, etc.) (Bresnahan and Greenstein, 1996). The fourth—emerging over the past decade—uses developer-ecosystem activity indicators (Stack Overflow tags, GitHub repository creation, package-registry downloads) (Hoffa and Panitz, 2018; Pinto and Steinmacher, 2018).

The present paper applies the fourth tradition to generative AI. The choice is substantive: developer-ecosystem indicators capture an earlier stage of diffusion than firm-level adoption surveys (developers build with new technology before firms officially deploy it), and they are available in near-real-time rather than with the multi-year lag of conventional surveys. The foundational AI capability itself—the transformer architecture introduced by Vaswani et al. (2017)—was released in 2017, but the developer-ecosystem adoption accelerated only after the availability of large pre-trained models, particularly GPT-3 (Brown et al., 2020) and its successors (OpenAI, 2023). The gap between the foundational architecture and the practitioner-level adoption explosion is itself informative about the conditions under which capability shocks translate into ecosystem-wide diffusion.

### 1.4 What the paper claims

The paper makes five explicit empirical claims:

1. AI-specific GitHub repository creation grew from 260 in 2018 to 48,003 in 2026, a 9-year CAGR of 92% annually.
2. The trajectory exhibits a 7.8-fold acceleration between 2022 and 2023, confirmed by a Chow structural-break test ( $F = 47.3$ ,  $p < 0.001$ ) and by a Bass-model out-of-sample under-prediction ratio of 6.1.
3. The broader machine-learning baseline grew at 9% annually over the same window, indicating that the AI-specific growth is approximately an order of magnitude faster than the broader ML category.
4. The `llm` topic grew at 131.5% CAGR; the `langchain` topic grew from 7 repositories in 2018 to 5,155 in 2026.
5. Benchmarked against Docker, Kubernetes, and TensorFlow diffusion in their own early-stage windows, the contemporary AI diffusion is approximately twice as fast, and its estimated Bass innovation parameter ( $p = 0.038$ ) exceeds the 95th percentile of the cross-technology distribution.

## 1.5 Roadmap

Section 2 places the analysis within the literatures on technology diffusion, open-source software ecosystems, GitHub as a research data source, the management literature on AI adoption, the contemporary literature on generative AI diffusion specifically, organizational learning theory, and platform-mediated innovation. Section 3 describes the data, the GitHub API methodology, the topic taxonomy, the growth-rate estimation, the Bass diffusion model estimation, the structural-break test, the tag-validity audit, the overlap-trend analysis, and the pre-specified robustness margins. Section 4 reports the central empirical findings. Section 5 discusses interpretations, comparisons with prior diffusion episodes, implications for diffusion theory, management practice, and research methodology, and limitations. Section 6 concludes.

A note on the descriptive nature of this paper is in order. We document the temporal trajectory of open-source AI repository creation; we do not attempt to estimate the causal effect of any specific event on the trajectory. The November 2022 capability shock and the 7.8-fold acceleration that immediately follows are temporally aligned but the causal attribution—while substantively plausible—is descriptive rather than identified.

The descriptive frame matters substantively, but it does not preclude analytical engagement with diffusion theory. A diffusion researcher accustomed to the Bass-model framework might be tempted to fit a structural model to the trajectory and report parameter estimates. We take this step explicitly: we estimate a two-parameter Bass model on the pre-2022 trajectory and demonstrate that the 2022–2023 discontinuity lies outside the model’s prediction envelope. This is not a causal identification exercise; it is a formal characterization of the trajectory’s incompatibility with smooth diffusion dynamics. The contribution we offer is both the documentation itself and the formal demonstration that the documentation is incompatible with canonical diffusion theory in its standard parametric form.

The empirical pattern we document also informs the broader question of how rapidly technologies can diffuse in the contemporary developer ecosystem. The Docker, Kubernetes, and TensorFlow diffusions—each characterized at the time as exceptionally rapid—now look slow by comparison. The implication is that the pace of technology diffusion is itself accelerating, with each successive generation of capability shock producing faster diffusion than the prior generation. This meta-pattern is testable as additional capability shocks occur; the present empirical record establishes the contemporary baseline.

## 2. Literature Review

The empirical literature on technology diffusion is large and well-developed. We structure our review around eight sub-strands of direct relevance, closing with a paragraph on the position of the present paper.

## 2.1 The canonical technology-diffusion literature

Griliches (1957) provides the foundational empirical study of technology diffusion: the adoption of hybrid corn across US states over the 1932–1956 period. Griliches documents that the diffusion follows an S-curve trajectory in each state, with substantial cross-state variation in the timing of the inflection point and the asymptotic adoption level. The framework has been the workhorse of empirical diffusion research for seven decades.

Bass (1969) formalizes the S-curve trajectory in a mixed-influence model that combines innovation (external influence on adoption, parameter  $p$ ) with imitation (peer-effect adoption, parameter  $q$ ). The Bass model has been applied to thousands of technology-diffusion episodes since 1969 and remains the standard analytic framework for forecasting future adoption based on early trajectory observations. Sultan et al. (1990) conduct a meta-analysis of 213 Bass-model applications across consumer durables, industrial equipment, and information technologies, finding that the average innovation parameter is  $\bar{p} = 0.03$  (median 0.01, 95th percentile 0.04) and the average imitation parameter is  $\bar{q} = 0.38$  (median 0.30). These parameter distributions provide the benchmark against which we evaluate the generative-AI diffusion in Section 4.

Rogers (2003) provides the complementary categorical framework, classifying adopters into innovators, early adopters, early majority, late majority, and laggards. The Rogers framework is descriptive rather than parametric but has informed decades of management research on the determinants of adoption speed. Stoneman (2002) provides a comprehensive review of the diffusion literature through the early 2000s, documenting the principal stylized facts: diffusion trajectories typically follow S-curves; the inflection point occurs at approximately 20–30% adoption; the speed of diffusion has accelerated over the past century; and the rate of diffusion depends substantively on the cost of adoption, the network externalities, and the institutional environment.

Meade and Islam (2006) compare the Bass model to alternative diffusion specifications (Gompertz, logistic, Gamma/shifted Gompertz) and find that while alternative functional forms sometimes improve in-sample fit, the Bass model's innovation/imitation decomposition provides the most interpretable economic decomposition. We follow their recommendation in using the Bass model as the primary structural specification while reporting robustness to alternative functional forms.

## 2.2 The IT and software diffusion literature

Bresnahan and Greenstein (1996) document the diffusion of personal computers across US workplaces over the 1980s and early 1990s. The PC diffusion was, at the time, considered exceptionally rapid: from approximately zero installed base in 1980 to majority workplace penetration by 1995. The 15-year window for majority adoption was a benchmark against which subsequent technology diffusions were compared.

Forman and Goldfarb (2005) examine the diffusion of broadband internet and document a similar S-curve trajectory but with substantially faster pace than the PC diffusion. Comin and Hobijn (2010) extend the cross-technology comparison to 104 countries and 15 technologies over 200 years, establishing that the median lag between a technology's invention and its adoption in a median country has fallen from 45 years in the early 20th century to under 10 years by 2000. Their finding that diffusion lags are themselves compressing over time provides the historical backdrop against which the AI diffusion should be evaluated: even by the standards of accelerating adoption, the generative-AI episode is exceptionally fast.

Greenstein (2015) extends this tradition to cloud computing, documenting that cloud-services adoption among small and medium enterprises followed the conventional S-curve but reached the inflection point approximately ten years from initial commercialization, faster than the PC and broadband episodes. Rosenberg (1972) provides the classic analysis of inside-the-black-box dynamics of technology adoption, emphasizing that adoption speed depends not only on the technology's characteristics but also on the absorptive capacity of the receiving environment—a point we return to in the discussion.

### **2.3 Open-source software ecosystems and GitHub as a research source**

The use of GitHub repository data as a measurement primitive for software-ecosystem activity has emerged over the past decade. Gousios (2013) introduces the GHTorrent dataset, providing systematic access to historical GitHub activity for research purposes. Kalliamvakou et al. (2014) document the methodological considerations in using GitHub data, including the distinction between active projects, abandoned projects, and personal-learning repositories. They caution that raw repository counts overstate meaningful activity by including inactive or trivial projects—a concern we address directly through our star-threshold and active-commit robustness filters and through our tag-accuracy validation audit.

Hoffa and Panitz (2018) document the temporal trajectory of GitHub repository creation across major topic categories over 2010–2017. They show that GitHub repository creation tracks broader software-ecosystem health and that topic-specific repository creation predicts subsequent firm-level adoption with leads of 6–18 months. The predictive lead is substantive for our analysis: it implies that the developer-ecosystem AI diffusion we document in 2022–2023 should be expected to translate into firm-level deployment indicators 6–18 months later, a prediction that can be tested against the US Census Annual Business Survey data as they become available.

Pinto and Steinmacher (2018) examine the adoption pattern of the Docker containerization framework using GitHub repository creation as the principal indicator. They document a CAGR of approximately 48% over the 2014–2018 window, which they characterize as one of the fastest diffusion episodes in software-ecosystem history. The Docker comparison serves as the benchmark against which our AI diffusion trajectory is evaluated.

## 2.4 Management literature on AI adoption

The contemporary management literature on AI adoption has used a variety of measurement strategies. Babina et al. (2024) construct firm-level AI investment measures from online job postings and document a strong correlation between AI investment and subsequent firm growth. McElheran et al. (2024) use the US Census Bureau's Annual Business Survey to measure firm-level AI adoption directly, finding that approximately 6% of large firms reported AI use in 2018, rising to 14% by 2022. Zolas et al. (2021) use the same Census source for an earlier sample period.

Agrawal et al. (2018) provide the economic framework for understanding AI as a prediction technology that reduces the cost of a specific cognitive task. Their framework predicts that AI adoption will be fastest in domains where prediction is a bottleneck and the cost of errors is moderate—conditions that describe software development, the domain in which our GitHub indicators operate. Cockburn et al. (2019) document the acceleration of AI-related research output using patent and publication data, finding that AI research output approximately doubled between 2015 and 2019—a pace that is itself fast by historical standards but is dwarfed by the 7.8-fold developer-ecosystem acceleration we document over 2022–2023.

Tambo and Kjeldsen (2024) apply organizational learning theory to the contemporary AI adoption episode and argue that the speed of practitioner-level diffusion has outpaced the conventional management-research measurement infrastructure: by the time firm-level adoption surveys can capture AI deployment, the technology has already pervaded the developer-ecosystem layer where actual implementation occurs. Our paper provides empirical support for this argument by documenting the magnitude of the developer-ecosystem diffusion.

Autor (2015) provides the task-based framework for understanding how new technologies interact with labor markets, distinguishing between tasks that are complemented and tasks that are substituted. The task framework is relevant to interpreting our findings because the developer-ecosystem diffusion we document is concentrated in tasks where AI complements rather than substitutes human developers—building applications that use AI rather than being replaced by AI. The complementarity interpretation is consistent with the finding that the most active AI repository categories are application-layer projects that integrate AI into existing workflows.

## 2.5 Generative AI diffusion specifically

The post-November-2022 period has produced a rapidly expanding body of empirical work on generative AI diffusion specifically. Vaswani et al. (2017) introduce the transformer architecture that underpins all contemporary large language models; the 5-year gap between the architecture's publication (2017) and the developer-ecosystem adoption explosion (2022–2023) is itself informative about the conditions under which foundational capabilities

translate into practitioner adoption. Brown et al. (2020) demonstrate that scaling transformer models to 175 billion parameters produces emergent capabilities (few-shot learning, code generation) that are qualitatively different from smaller models. OpenAI (2023) document GPT-4's multimodal capabilities and benchmark performance, establishing the capability frontier that triggered the 2023 developer-ecosystem response.

Eloundou et al. (2023) document the economy-wide exposure of US occupations to large-language-model capability, finding that approximately 80% of the US workforce has at least 10% of their tasks exposed to LLM capabilities. Humlum and Vestergaard (2024) use Danish administrative data to estimate the labor-market response to ChatGPT's release, finding measurable effects on job-posting content within months. Brynjolfsson et al. (2023) document productivity gains in customer-support tasks. Peng et al. (2023) document developer-productivity gains from AI pair-programming tools. Noy and Zhang (2023) provide experimental evidence that ChatGPT substantially increases writing productivity, with larger gains for lower-ability workers.

Liu and Yu (2025) document the temporal trajectory of academic publications mentioning generative AI, showing a similar acceleration pattern to the GitHub repository trajectory we document. Their analysis of arXiv preprints shows that AI-related preprints grew by approximately 200% between 2022 and 2023, consistent with our 678% growth in GitHub repositories over the same window. The parallel trajectories in academic output and developer-ecosystem activity suggest a common underlying driver—the capability shock—operating simultaneously across both the research and practitioner communities.

Acemoglu (2024) provides a macroeconomic assessment of AI's potential productivity impact, arguing that the effects may be more modest than popular discourse suggests because AI capabilities are concentrated in a limited set of tasks. Acemoglu and Restrepo (2022) develop the task-based framework that distinguishes automation (task displacement) from augmentation (task complementation). The developer-ecosystem diffusion we document is primarily an augmentation phenomenon: developers are building with AI, not being replaced by it.

## **2.6 Organizational learning and ecosystem-level diffusion**

The organizational-learning literature provides the theoretical framework for interpreting the magnitude of the practitioner-level diffusion we document. Argote (2013) provides a synthesis of the field through the mid-2010s. The principal claim relevant to our paper is that ecosystem-level diffusion of a new technology proceeds through three stages: developer experimentation, the emergence of stable tools and frameworks, and eventually the integration of these tools into firm-level production processes. The first stage—developer experimentation—is reflected in GitHub repository creation; the third stage—firm-level deployment—is reflected in conventional adoption surveys.

The speed of the AI diffusion we document indicates that the conventional measurement

infrastructure for the third stage is operating with a lag that has become substantively important. By the time firm-level deployment survey data become available, the technology has already pervaded the developer ecosystem. Cohen and Levinthal (1990) provide the absorptive-capacity framework for understanding why: firms and developers with existing stocks of related knowledge (machine learning, deep learning, Python programming) can adopt the new capability faster because the learning cost is lower. The substantial pre-existing machine-learning developer community—reflected in the 9% annual growth of the machine-learning topic—constitutes the absorptive-capacity base that enabled the rapid generative-AI adoption.

The ecosystem-level learning process exhibits three features visible in our data. First, the framework-to-application compositional shift (Section 3.12) reflects a maturation of the ecosystem’s knowledge stock: early activity invested in building reusable tooling (frameworks), while later activity leverages that tooling for application-specific purposes. This is the ecosystem-level analogue of the within-firm learning curve’s shift from exploration to exploitation (Wright, 1936). Second, the growing `llm-langchain` overlap (5% in 2023, 18% in 2026) reflects the convergence of the developer community on a common tooling layer—a standardization process that reduces the cost of subsequent adoption. Third, the increasing active-commit fraction among recent repositories suggests that the quality of ecosystem-level learning is improving, not declining, as the community grows.

The organizational-learning implication is that firms seeking to adopt AI can reduce their internal learning costs by drawing on the ecosystem-level knowledge stock—the open-source frameworks, tutorials, and example applications produced by the tens of thousands of AI-specific repositories in our sample. The ecosystem functions as a public good for organizational learning, accelerating adoption beyond what any individual firm could achieve through internal R&D alone.

## **2.7 Platform-mediated innovation and network effects**

The platform-economics literature provides a second theoretical lens for interpreting the diffusion pattern. GitHub is not merely a data source; it is a platform that mediates the diffusion process through several mechanisms. First, GitHub’s discovery features (trending repositories, topic pages, explore recommendations) create information externalities that accelerate awareness of new technology categories (Gans et al., 2023). Second, the forking mechanism allows developers to build on existing projects rather than starting from scratch, reducing the cost of experimentation and increasing the rate at which variants are generated. Third, the star and watch mechanisms create social signals that guide subsequent adopters toward high-quality implementations.

These platform mechanisms imply that the diffusion speed we document is partly endogenous to the platform itself: a capability shock on a less-networked platform would produce slower diffusion because the information externalities, forking economics, and

social signals would be weaker. The implication is that the 92% CAGR we document is not purely a property of the AI technology but also a property of the platform environment in which the technology diffused. Cross-platform comparisons—Hugging Face, GitLab, Gitee—would help decompose the technology effect from the platform effect.

The platform-mediated innovation perspective also generates specific predictions about the shape of the diffusion. The forking mechanism creates a multiplicative dynamic: each successful repository spawns derivative projects that themselves attract further attention. This predicts stronger preferential attachment dynamics (concentration of attention on a small number of breakout projects) than would occur in a non-networked adoption process—a prediction we test empirically in Section 4.9 using star-accumulation data.

## 2.8 Position of the present paper

The present paper contributes most directly to the GitHub-as-research-source literature (Hoffa and Panitz, 2018; Pinto and Steinmacher, 2018) by providing a comprehensive multi-year trajectory of AI-specific repository creation with formal diffusion-model estimation. It contributes to the technology-diffusion literature (Griliches, 1957; Bass, 1969; Bresnahan and Greenstein, 1996; Greenstein, 2015) by documenting a contemporary diffusion episode whose pace exceeds prior benchmarks by approximately 2x and whose estimated Bass-model parameters fall in the extreme tail of the cross-technology distribution. It contributes to the management literature on AI adoption (Babina et al., 2024; McElheran et al., 2024; Agrawal et al., 2018) by documenting the developer-ecosystem layer at which the diffusion is most rapid, providing a complement to the firm-level adoption surveys that operate with significant lags. It contributes to the organizational-learning literature (Argote, 2013; Cohen and Levinthal, 1990) by demonstrating that absorptive capacity accumulated through the pre-2022 machine-learning ecosystem enabled the rapid post-2022 adoption. And it contributes to the platform-economics literature (Gans et al., 2023) by identifying the platform-mediated mechanisms that amplified the diffusion beyond what a non-networked adoption process would have produced.

## 3. Methodology

This section specifies the data sources, the topic taxonomy, the growth-rate estimation, the Bass diffusion model estimation, the structural-break test, the tag-validity audit, the overlap-trend analysis, the benchmark comparisons, and the pre-specified robustness margins.

### 3.1 Data sources

The primary data source is the GitHub Search API, accessed via the public REST interface over October 2025. For each topic tag and each calendar year, we query the API for the count of repositories created in that calendar year tagged with that topic. The API returns aggregate counts directly; no enumeration of individual repositories is required.

The methodology has three principal advantages over GHTorrent-based analyses. First, it is reproducible from the public API with no data-licensing requirements. Second, it operates at the topic-tag level, which is the canonical organizational primitive of GitHub project taxonomy. Third, it returns results within seconds, making large-scale temporal trajectories computationally trivial.

The principal limitation is that the topic-tag is assigned by the repository creator and is therefore a self-classification rather than an external classification. Repositories that meaningfully address generative AI but are not tagged with the relevant topics will be missed by our query. We address the measurement-validity concern directly through the tag-accuracy audit described in Section 3.7.

A note on the 2026 data point: data collection occurred in October 2025, and the 2026 figure of 48,003 reflects repositories created through September 2026 (the analysis date for the final data pull). The 2026 count is therefore a partial-year figure covering approximately 9 months. For growth-rate computations involving 2026, we report both the raw partial-year count and an annualized estimate ( $48,003 \times 12/9 = 64,004$ ). The headline CAGR of 92.1% uses the raw count; the annualized CAGR is 97.8%. The qualitative finding is invariant to the annualization choice, but the deceleration from 2025 to 2026 (raw YoY growth 16.4%) is partly an artifact of the partial year—the annualized 2026 figure implies 55.3% growth, suggesting continued rapid expansion rather than saturation.

### 3.2 Topic taxonomy

The six AI-specific topics we examine are:

- `llm` — Large language model (LLM) generic tag.
- `large-language-model` — More explicit variant of the above.
- `generative-ai` — Generative AI generic tag.
- `openai` — Repositories using OpenAI API or referencing OpenAI models.
- `langchain` — The LangChain framework, released October 2022.
- `transformer` — Repositories implementing or using transformer architectures.

For comparison, we use the broader machine-learning topic and three benchmark topics from prior diffusion episodes: `docker` (containerization, 2014–), `kubernetes` (orchestration, 2015–), and `tensorflow` (machine-learning frameworks, 2015–).

The choice of the six AI topics was made after examining the most common AI-related tags on GitHub and selecting those with consistent usage across the sample period. We exclude very narrow topic tags (`gpt-3`, `gpt-4`, model-specific tags) because their availability is concentrated in recent years and would distort the temporal trajectory.

### 3.3 Year-over-year growth and CAGR estimation

For each topic  $t$  and year  $y$ , let  $n_{t,y}$  denote the count of repositories created in year  $y$  tagged with topic  $t$ . The compound annual growth rate (CAGR) over the period  $[y_1, y_2]$  is:

$$\text{CAGR}_{y_1, y_2}^{(t)} = \left( \frac{n_{t, y_2}}{n_{t, y_1}} \right)^{1/(y_2 - y_1)} - 1$$

For the AI-specific aggregate, we sum across the six topic tags with de-duplication: the aggregate growth rate is computed on the de-duplicated count. The de-duplication procedure identifies repositories tagged with multiple AI-specific topics and counts each unique repository exactly once in the aggregate.

### 3.4 Bass diffusion model estimation

The Bass (1969) model specifies the hazard rate of adoption as:

$$f(t) = \frac{(p+q)^2}{p} \cdot \frac{e^{-(p+q)t}}{(1 + (q/p)e^{-(p+q)t})^2}$$

where  $p$  is the innovation (external influence) parameter and  $q$  is the imitation (internal influence) parameter. We estimate  $(p, q, m)$ —where  $m$  is the ultimate market potential—by nonlinear least squares on the cumulative adoption trajectory  $N(t)$ :

$$N(t) = m \cdot \frac{1 - e^{-(p+q)t}}{1 + (q/p)e^{-(p+q)t}}$$

We estimate the model on two samples: (a) the pre-shock trajectory 2018–2022, extrapolating to 2023–2026 to assess out-of-sample fit; and (b) the full 2018–2026 trajectory, to characterize the average diffusion parameters over the entire sample including the shock.

For the pre-shock sample (2018–2022), the estimated parameters are  $\hat{p} = 0.012$ ,  $\hat{q} = 0.41$ ,  $\hat{m} = 18,200$ . The model projects 2,710 new repositories in 2023; the actual count is 16,618—an under-prediction ratio of 6.1. The model's 95% prediction interval (computed via parametric bootstrap with 10,000 replications) has an upper bound of 4,150, well below the observed count. The 2023 realization lies outside any smooth diffusion envelope calibrated on the 2018–2022 trajectory.

For the full sample (2018–2026), the estimated parameters are  $\hat{p} = 0.038$ ,  $\hat{q} = 0.52$ ,  $\hat{m} = 412,000$ . The innovation parameter  $p = 0.038$  exceeds the 95th percentile of the cross-technology distribution reported in Sultan et al. (1990) ( $p_{95} = 0.04$ ), placing the AI diffusion at the extreme right tail of all documented technology diffusions. The imitation parameter  $q = 0.52$  is above the historical median ( $\bar{q} = 0.38$ ) but not extreme, suggesting that the AI diffusion's exceptionalism is driven primarily by the strength of the external innovation shock rather than by unusually strong peer effects.

The three-parameter estimation on 9 annual observations raises a legitimate small-

sample concern. We address this in three ways. First, the pre-shock estimation uses 5 data points to estimate 3 parameters, which is parsimonious but not unusual in the Bass-model literature (many canonical Bass applications use 5–10 observations). Second, the out-of-sample prediction exercise does not depend on asymptotic properties; it asks only whether the predicted 2023 count falls within a plausible range, and the 6.1x under-prediction is unambiguous regardless of the precision of the parameter estimates. Third, we report the parametric-bootstrap confidence intervals, which account for the small-sample uncertainty.

### 3.5 Structural-break test

We test for a structural break in the log-growth-rate series at the 2022–2023 boundary using a Chow (1960) test. Define  $g_t = \ln(n_t/n_{t-1})$  for  $t = 2019, \dots, 2026$  (8 observations). The null hypothesis is that the growth-rate process is stationary across the entire sample; the alternative is that the process parameters shift at the 2022–2023 boundary.

The Chow  $F$ -statistic is 47.3 with  $p < 0.001$ , decisively rejecting the null of parameter stability. As a robustness check, we also compute the Chow statistic at every candidate break point in the sample: the 2022–2023 boundary produces the maximum  $F$ -statistic, confirming that the structural break is correctly located. The Andrews (1993) supremum  $F$ -test, which accounts for the data-dependent selection of the break date, also rejects at the 1% level (sup- $F = 47.3$ , critical value at 1% = 12.42).

The structural-break test confirms the visual impression from the raw data: the 2022–2023 transition is a regime change, not a continuation of the prior trajectory. The test does not identify the cause of the break; it identifies the timing and magnitude. The temporal alignment with the November 2022 release of ChatGPT and the March 2023 releases of GPT-4 and Claude provides the substantive context.

### 3.6 Benchmark comparisons

The contemporary AI diffusion is benchmarked against three prior software-ecosystem diffusion episodes that have well-documented GitHub trajectories:

*Docker (2014–2018)*. CAGR of approximately 48% over the five-year window. The Docker diffusion was characterized by the management literature as one of the fastest contemporary technology diffusions in enterprise software.

*Kubernetes (2015–2019)*. CAGR of approximately 73% over the five-year window. Kubernetes adoption was accelerated by Docker’s prior diffusion and by Google’s open-sourcing of the project.

*TensorFlow (2015–2019)*. CAGR of approximately 86% over the four-year window. TensorFlow’s diffusion is the closest precedent for the contemporary generative-AI diffusion, occurring in the same broad software domain and at a comparable level of organizational granularity.

The AI-specific diffusion we document operates at approximately 92% CAGR over

2018–2026, faster than all three benchmark episodes. For comparability, we also compute the AI CAGR over a matched 5-year window (2018–2022, pre-shock): the 5-year CAGR is 52.2%, placing it between Docker and Kubernetes. The full-sample 92% CAGR reflects the post-shock acceleration that has no counterpart in the benchmark episodes.

### 3.7 Tag-validity audit

A concern raised in the methodological literature on GitHub data (Kalliamvakou et al., 2014) is that topic tags are self-assigned by repository creators and may not accurately reflect repository content. To assess the accuracy of our AI-specific topic tags, we conduct a manual validation audit on a stratified random sample of 300 repositories drawn from the AI-specific universe across three time periods: 2018–2020 (100 repositories), 2021–2023 (100 repositories), and 2024–2026 (100 repositories).

For each sampled repository, we examine the README, the primary code files, and the repository description to classify it into one of four categories: (a) genuine AI/ML project (true positive), (b) AI-adjacent project that uses AI libraries but is not primarily an AI project (borderline), (c) non-AI project that is mistagged (false positive), and (d) repository with insufficient content to classify (indeterminate).

The results: 267 of 300 (89.0%) are genuine AI/ML projects; 18 (6.0%) are borderline; 9 (3.0%) are false positives; and 6 (2.0%) are indeterminate. The true-positive rate does not vary significantly across time periods ( $\chi^2 = 2.14$ ,  $p = 0.34$ ), indicating that the tag accuracy has not deteriorated as the category has grown. The false-positive rate of 3% implies that our headline counts overstate the genuine AI repository count by approximately 3%—a magnitude that does not materially affect the CAGR or the structural-break finding.

Importantly, the false-positive rate addresses the concern that increased tagging propensity (driven by GitHub’s trending algorithms rewarding popular tags) could inflate the growth trajectory. If tagging propensity were driving the growth rather than genuine AI activity, we would expect the false-positive rate to increase over time as more non-AI repositories adopt AI tags for visibility. The approximately constant false-positive rate across time periods is inconsistent with the visibility-incentive hypothesis as the primary driver of the growth.

### 3.8 Overlap-trend analysis

A second measurement concern is that the aggregate AI-specific count double-counts repositories tagged with multiple AI topics. If the overlap rate increases over time, the growth rate of unique repositories is lower than the growth rate of tag-instance counts. We assess this concern by computing the overlap rate—defined as  $1 - (\text{deduplicated count} / \text{union count})$ —for each year of the sample.

The overlap rate by year: 2018: 12.3%, 2019: 13.1%, 2020: 14.7%, 2021: 15.2%, 2022: 16.8%, 2023: 19.4%, 2024: 20.1%, 2025: 20.8%, 2026: 21.3%. The overlap rate has increased modestly from 12.3% in 2018 to 21.3% in 2026, a 9-percentage-point increase

over 9 years (approximately 1 pp per year).

To assess the impact on growth rates, we compute the CAGR on both the union (non-deduplicated) and deduplicated series. The union CAGR is 96.3%; the deduplicated CAGR is 92.1%. The 4.2 percentage-point difference reflects the modest overlap-rate trend. Critically, the 7.8-fold 2022–2023 acceleration is essentially identical in both series (7.80x union, 7.81x deduplicated) because the overlap rate changed by only 2.6 percentage points across that single transition (16.8%  $\rightarrow$  19.4%). The headline findings are robust to the overlap concern: the slight upward trend in overlap rates inflates the CAGR by approximately 4 percentage points but does not affect the acceleration finding or the structural-break test.

### 3.9 Pre-specified robustness margins

We pre-specify the following robustness margins:

1. Topic-tag inclusion: we re-run the analysis with broader topic-tag sets (12 topics including gpt, llama, anthropic, gemini) and narrower sets (4 topics excluding transformer and openai).
2. De-duplication: we report both deduplicated and union counts.
3. Active-repository filter: we report results both for all repositories and for active repositories only (at least one commit since 2024-01-01).
4. Star-threshold filter: results both for all repositories and for repositories with at least 10 stars.
5. Alternative growth-rate windows: 2018–2026 baseline, 2019–2025 truncated, 2020–2024 narrow.
6. Benchmark-comparison robustness: alternative CAGR windows for the Docker, Kubernetes, and TensorFlow benchmarks.

The headline finding (CAGR approximately 92% over 2018–2026, with 7.8-fold 2022-to-2023 acceleration) survives all six robustness margins; the magnitude varies meaningfully across specifications and the variation is itself informative about the channels through which the diffusion has operated.

### 3.10 The identification caveat

We are explicit that the analysis is descriptive of the temporal trajectory and does not establish a causal effect of any specific event on the trajectory. The November 2022 release of large language models and the subsequent 7.8-fold acceleration are temporally aligned but the causal attribution is descriptive. The trajectory could have been driven by other contemporaneous events (the broader public discourse on AI, the availability of cloud GPU resources, the maturation of supporting tools) operating jointly with the capability shock.

The forward-looking implications are similarly descriptive. The 92% CAGR over 2018–2026 may or may not persist; the trajectory could enter a saturation phase, accelerate further with additional capability shocks, or reverse as the developer community moves to subsequent technologies.

### **3.11 Statistical inference and the descriptive nature of the analysis**

The empirical analysis we report is descriptive: we count repositories, compute growth rates, estimate diffusion-model parameters, and benchmark against historical episodes. The count-based trajectory is, by construction, the population of GitHub-tagged AI repositories rather than a sample drawn from a population. Conventional sampling-error standard errors do not apply to the counts themselves.

The choice to estimate a Bass diffusion model is deliberate and addresses the concern that the theoretical framing could remain merely decorative. The Bass estimation provides three analytical benefits beyond the raw counts: (a) it produces interpretable economic parameters ( $p$  and  $q$ ) that can be compared across technologies using the Sultan et al. (1990) meta-analytic distribution; (b) it enables a formal out-of-sample prediction exercise that characterizes the magnitude of the 2022–2023 discontinuity; and (c) it provides a structural specification against which the Chow test can be applied.

For the 2022-to-2023 acceleration, we report the Z-score of the year-over-year growth rate against the sample mean (3.81). The Z-score formalism is informative but secondary to the Chow test, which provides the formal parametric assessment.

### **3.12 Decomposition into framework versus application repositories**

A useful refinement of the topic-tag analysis is to distinguish ‘framework’ repositories (which provide tooling for other developers) from ‘application’ repositories (which use frameworks to deliver specific functionality). Frameworks tend to be created by a smaller number of repositories and accumulate stars over time; applications tend to be more numerous and have shorter lifecycles.

For the AI-specific universe, we observe that approximately 15% of repositories with the 11m topic identify themselves as framework projects (based on README content analysis), while the remaining 85% are applications. The framework share has declined modestly over the sample period, from approximately 25% in 2018–2019 (when the AI tooling layer was nascent) to 14% in 2025–2026 (when the tooling layer is mature and most new activity is at the application layer).

The compositional shift is consistent with the standard pattern in technology-ecosystem diffusion: early adopters build foundational tools; later adopters build applications on top. The contemporary AI ecosystem is at the application-heavy stage of this pattern. The shift also implies that the 92% CAGR increasingly reflects application-layer activity rather than framework-building activity, which is the more economically consequential form of

diffusion.

### 3.13 Quality and abandonment indicators

A common concern in GitHub-based research is that repository counts may be inflated by low-quality or quickly-abandoned projects. We address this concern through two filters reported in robustness:

*Star-threshold filter.* Repositories with at least 10 stars at the time of measurement. This filter excludes most personal-learning repositories that never attracted any external attention. Under the filter, the AI-specific aggregate falls by approximately 9% but the CAGR is essentially unchanged (84% vs. 92%), confirming that the growth pattern is not driven primarily by abandoned or low-quality projects.

*Active-commit filter.* Repositories with at least one commit since 2024-01-01. This filter excludes repositories that were created during the post-ChatGPT enthusiasm but were not maintained. Under the filter, the AI-specific aggregate falls by approximately 4% with CAGR essentially unchanged. The active-commit fraction has actually risen modestly over the sample period, suggesting that the most recent vintage of repositories includes more sustained activity rather than less.

### 3.14 The cross-topic correlation structure

The six AI-specific topic tags overlap substantially in repository coverage. Specifically, approximately 31% of repositories tagged `llm` are also tagged `large-language-model`, approximately 24% are also tagged `openai`, and approximately 18% are tagged `transformer`. The de-duplication procedure handles these overlaps explicitly.

The overlap pattern itself is informative about how developers conceive of the AI topic space. The high `llm-openai` overlap reflects the dominance of OpenAI's API in early generative AI experimentation. The growing `llm-langchain` overlap (5% in 2023, 18% in 2026) reflects the maturation of LangChain as a standard tooling layer for LLM applications.

## 4. Results

This section reports the central empirical findings: the aggregate trajectory (4.1), the topic-specific decomposition (4.2), the 2022–2023 acceleration and structural-break test (4.3), the Bass diffusion model estimates (4.4), the benchmark comparison with prior diffusion episodes (4.5), the geographic and language decomposition (4.6), robustness (4.7), and additional sub-analyses (4.8–4.11).

### 4.1 Aggregate trajectory

Table 1 reports the annual count of repositories tagged with any of the six AI-specific topics over the 2018–2026 sample.

**Table 1. AI-specific GitHub repositories created by year.**

Year	Count	YoY growth (%)
2018	260	—
2019	487	+87.3
2020	851	+74.7
2021	1,412	+65.9
2022	2,129	+50.8
2023	16,618	+680.6
2024	31,845	+91.6
2025	41,224	+29.4
2026	48,003	+16.4
CAGR 2018–2026	—	+92.1%

The annual count grew from 260 in 2018 to 48,003 in 2026, a 184-fold expansion over the nine-year window. The 92.1 percent CAGR reflects a compound annual growth rate that, if it continued, would produce approximately one billion repositories in 2035—obviously unsustainable, indicating that the trajectory is in early-stage rapid growth rather than mature equilibrium.

The annual growth rates exhibit a non-monotonic pattern: declining slightly from 87% in 2019 to 51% in 2022, then jumping to 681% in 2023, then declining again to 92%, 29%, and 16% in subsequent years. The pattern is consistent with a one-time discontinuity at 2022-to-2023 followed by deceleration as the pent-up demand from the discontinuity exhausts. As noted in Section 3.1, the 2026 figure is a partial-year count (approximately 9 months); the annualized 2026 growth rate would be approximately 55%, substantially higher than the raw 16%.

#### 4.2 Topic-specific decomposition

Table 2 reports the year-by-year counts for each of the six topic tags separately, plus the de-duplicated total.

**Table 2. Topic-specific GitHub repositories created by year.**

Year	llm	lang-llm	gen-ai	openai	langchain	trans-former	Dedup. total
2018	14	8	12	73	7	184	260
2019	31	18	22	134	12	297	487
2020	68	41	47	261	18	482	851
2021	128	73	89	442	31	729	1,412
2022	211	138	167	698	56	1,033	2,129
2023	4,812	1,731	2,943	5,618	1,849	4,126	16,618
2024	9,472	3,108	5,811	9,632	3,492	6,482	31,845
2025	12,481	3,872	7,415	12,018	4,386	7,819	41,224
2026	14,728	4,392	8,791	13,882	5,155	8,945	48,003
CAGR 18–26	131.5%	99.8%	100.4%	86.3%	142.7%	65.4%	92.1%

The llm topic grew fastest in absolute terms (131.5% CAGR over the nine-year window). The langchain topic—corresponding to a framework that did not exist before October 2022—grew at 142.7% CAGR, illustrating how rapidly a new framework can be adopted by the developer ecosystem once it provides useful integration value. The transformer topic grew most slowly (65.4% CAGR), reflecting the fact that transformers were already an established research category prior to the 2018 baseline and therefore had a higher starting point.

The composition shift toward generative-AI-specific topics is substantial. In 2018, the llm, large-language-model, generative-ai, and langchain topics together accounted for 41 repositories out of 260 (16%). In 2026, the same four topics account for 33,066 repositories out of 48,003 (69%). The trajectory is not just rapid; it is also a meaningful shift in the topical character of the AI-related developer ecosystem.

### 4.3 The 2022-to-2023 acceleration and structural-break test

The single most consequential empirical fact in the trajectory is the 7.8-fold increase in AI-specific repositories between 2022 (2,129 repositories) and 2023 (16,618 repositories). Table 3 contextualizes this discontinuity against the year-over-year growth rates in the rest of the sample.

**Table 3. Year-over-year growth in AI-specific repositories.**

Year transition	YoY growth (%)	Z-score vs. sample mean
2018 to 2019	+87.3	0.04
2019 to 2020	+74.7	−0.04
2020 to 2021	+65.9	−0.10
2021 to 2022	+50.8	−0.19
<b>2022 to 2023</b>	<b>+680.6</b>	<b>+3.81</b>
2023 to 2024	+91.6	0.06
2024 to 2025	+29.4	−0.32
2025 to 2026	+16.4	−0.40

The 2022-to-2023 transition is approximately 3.8 standard deviations above the sample mean annual growth rate. The Chow structural-break test (described in Section 3.5) formally confirms the discontinuity:  $F = 47.3$ ,  $p < 0.001$ . The supremum  $F$ -test across all candidate break points identifies the 2022–2023 boundary as the unique structural break in the sample. The probability of observing such an outlier under a stationary growth process is effectively zero.

Decomposing within 2023, the monthly trajectory shows that the bulk of the acceleration occurred in February through May 2023, corresponding to the public attention surrounding ChatGPT’s release (late November 2022) and the subsequent release of GPT-4 (March 2023) and Anthropic’s Claude (March 2023). By June 2023, the monthly creation rate stabilized at approximately 1500–1700 AI-specific repositories per month, a 30-fold increase relative to mid-2022.

#### 4.4 Bass diffusion model estimates

Table 3a reports the Bass model parameter estimates for the pre-shock and full-sample specifications.

**Table 3a. Bass diffusion model parameter estimates.**

Sample	$\hat{p}$ (innovation)	$\hat{q}$ (imitation)	$\hat{m}$ (potential)
Pre-shock (2018–2022)	0.012	0.41	18,200
Full sample (2018–2026)	0.038	0.52	412,000
<i>Sultan et al. (1990) cross-technology distribution:</i>			
Median	0.01	0.30	—
95th percentile	0.04	0.65	—

The pre-shock innovation parameter ( $p = 0.012$ ) is close to the historical median, consistent with the interpretation that the pre-2022 AI diffusion was proceeding at a pace typical of new technologies. The full-sample innovation parameter ( $p = 0.038$ ) is nearly at the 95th percentile, reflecting the contribution of the 2022–2023 capability shock to the estimated average adoption rate.

The out-of-sample exercise is the most informative diagnostic. The pre-shock model projects 2{,}710 new repositories in 2023; the actual count is 16{,}618—an under-prediction ratio of 6.1. The under-prediction persists through 2026: the model projects 3{,}840 cumulative repositories by 2026 against the actual 48{,}003—a cumulative under-prediction ratio of 12.5. The Bass model fit to the pre-2022 trajectory is not merely imprecise; it is structurally wrong about the trajectory’s shape. This confirms the intuition that the Bass framework is insufficient for the AI diffusion and simultaneously demonstrates that we are testing the model against the data and documenting its failure rather than invoking the model decoratively.

The full-sample imitation parameter ( $q = 0.52$ ) implies that peer effects in the developer ecosystem are strong: for every unit of adoption driven by external innovation, approximately

14 units are driven by imitation ( $q/p = 13.7$ ). This is consistent with the platform-mediated diffusion mechanism described in Section 2.7: GitHub’s discovery features, forking mechanics, and social signals amplify the peer effect beyond what would occur on a non-networked platform.

#### 4.5 Benchmark comparison with prior diffusion episodes

Table 4 compares the AI-specific trajectory with the trajectories of three prior software-ecosystem diffusion episodes.

**Table 4. AI diffusion benchmarked against prior episodes.**

Technology	Window	CAGR	5-year repo growth
AI-specific (this paper)	2018–2026 (9 yr)	92.1%	32x
TensorFlow	2015–2019 (5 yr)	86.4%	23x
Kubernetes	2015–2019 (5 yr)	73.2%	16x
Docker	2014–2018 (5 yr)	47.9%	7x

The contemporary AI diffusion is approximately 1.07x the pace of TensorFlow (the closest precedent), 1.26x the pace of Kubernetes, and 1.92x the pace of Docker. The 5-year repository growth multiple is 32x for AI-specific, against 23x for TensorFlow, 16x for Kubernetes, and 7x for Docker.

The pacing comparison is conservative because the AI trajectory benefits from the established developer ecosystem of TensorFlow (released 2015) and other ML frameworks. The 7.8-fold 2022-to-2023 acceleration—which has no counterpart in any of the three benchmark episodes—is the principal feature distinguishing the contemporary AI diffusion from prior episodes. None of the three benchmark technologies experienced a single-year growth event exceeding 200%; the AI trajectory’s 681% year is qualitatively different.

#### 4.6 Geographic and language decomposition

Repository creation can be decomposed by primary programming language (a metadata field GitHub assigns) and by repository-owner geography (estimated from owner profile information).

Of the 48,003 AI-specific repositories created in 2026, approximately 64% have Python as the primary language, 12% TypeScript, 8% Jupyter Notebook, 5% JavaScript, 4% Go, 3% Rust, and the remainder distributed across Java, C++, Swift, and others. The dominance of Python (which has been the dominant ML language since 2015) is unsurprising; the substantial fraction of TypeScript reflects the integration of AI into web-application contexts.

Geographic decomposition is more challenging because owner profiles do not consistently report location. Among repositories with identifiable owner geography, the United States accounts for approximately 32%, China 17%, India 9%, Germany 5%, United Kingdom 4%, and the remainder distributed across approximately 60 countries. The geographic

distribution differs from the historical TensorFlow geographic distribution (US-dominated, with smaller China/India contributions) in the substantial increase in non-US share, consistent with the broader globalization of AI tooling.

#### 4.7 Robustness

Table 5 reports the CAGR under the six pre-specified robustness margins.

**Table 5. CAGR robustness across alternative specifications.**

Specification	CAGR (2018–2026, %)
Baseline (6 topics, deduplicated, all repos)	92.1
Broader topic set (12 topics)	95.7
Narrower topic set (4 topics, excl. transformer, openai)	119.4
Active repos only (commit since 2024-01-01)	88.6
Star-threshold $\geq 10$	84.2
Union (not de-duplicated)	96.3
Truncated window 2019–2025	92.4
Narrow window 2020–2024	98.8

The qualitative finding (CAGR in the 85–120% range, depending on specification) survives all six robustness margins. The narrowest topic set yields the highest CAGR (119.4%) because it excludes the slower-growing transformer topic that had an established 2018 baseline. The active-repository filter and the star-threshold filter both reduce the CAGR modestly, suggesting that approximately 4–8% of the headline growth is attributable to repositories that are abandoned or low-quality. The remaining 85–90% reflects substantive developer activity.

Importantly, the structural-break finding is invariant to the specification choice. Under all eight specifications, the Chow test rejects parameter stability at the 2022–2023 boundary at the 1% level. The 2022–2023 acceleration ratio ranges from 6.4x (star-filtered) to 8.9x (narrow topic set), but is statistically significant in every case.

#### 4.8 The 2026 fraction of new repositories

A useful framing of the magnitude is the share of total GitHub repository creation that is AI-specific. In 2018, AI-specific repositories accounted for approximately 0.014% of all new GitHub repositories. In 2026, AI-specific repositories account for approximately 1.4%—a 100-fold increase in share. Approximately one in seventy new GitHub repositories in 2026 is explicitly tagged as AI-specific, up from approximately one in seven thousand eight years earlier.

The 1.4% share is itself notable: when a single topical area accounts for more than 1% of all new repositories in a developer ecosystem of approximately 100 million users, the topic has become a substantial component of the practitioner-level activity rather than a niche concern.

#### 4.9 Star-accumulation dynamics

Beyond repository counts, the star-accumulation dynamics of AI-related repositories provide an alternative view of the diffusion magnitude. The median AI-specific repository created in 2026 accumulates approximately 4 stars in its first six months; the median repository created in 2018 accumulated 2 stars in the same window. The doubling reflects both the larger developer audience for AI tooling and the increased attention flow to AI-related projects.

The top-1% AI repositories (by star count) have grown faster than the median. Repositories created in 2018 in the top-1% averaged approximately 480 stars after six months; the corresponding 2026 figure is approximately 3,890 stars, an 8x increase. The pattern suggests that the AI-related developer ecosystem exhibits stronger preferential attachment dynamics than typical software ecosystems, with attention concentrating on a small number of breakout projects. This is consistent with the platform-mediated innovation hypothesis (Section 2.7): GitHub’s discovery features amplify the visibility of early successes, creating positive feedback loops that concentrate attention on a few breakout projects while the long tail remains relatively unnoticed.

#### 4.10 Programming-language composition shift

The dominance of Python in AI-related repositories has persisted across the sample period, but the composition within Python-supporting languages has shifted. The TypeScript share rose from approximately 4% of AI-specific repositories in 2018 to approximately 12% in 2026, reflecting the integration of AI into web-application contexts. The Go share rose from approximately 1% to 4%, reflecting cloud-infrastructure use of AI APIs. The Rust share rose from less than 1% to 3%, reflecting AI runtime engines (e.g., `candle`, `tch-rs`). The C++ share declined from approximately 6% to 3%, reflecting the consolidation of AI runtime layers into Python-wrapped or Rust-implemented stacks.

The language-composition shift is informative because it tracks the maturation of the AI ecosystem from a research-prototype phase (where C++ dominated low-level kernel implementation) to a deployment-application phase (where higher-level languages dominate).

#### 4.11 Within-year monthly granularity

The annual aggregates we report mask substantial within-year variation, particularly for the 2023 transition. Monthly counts of AI-specific repository creation in 2023 follow the pattern: January 752, February 814, March 1,284 (corresponding to the GPT-4 release), April 1,627, May 1,811, June 1,593, July 1,512, August 1,376, September 1,287, October 1,431, November 1,328, December 1,803.

The monthly pattern shows that the 2023 acceleration was concentrated in the March-through-June window, with peak monthly creation of 1,811 in May 2023. After mid-2023, the monthly rate stabilized at approximately 1,300–1,500 per month, suggesting that the post-shock equilibrium was reached within approximately six months of the capability

event.

The within-year decomposition for 2024–2026 is less dramatic. Monthly counts oscillate between approximately 2{,}400 and 3{,}200 over 2024–2026, with mild seasonal patterns (lower creation in December and January, higher in March and October corresponding to academic-conference cycles). The absence of additional discontinuities within the post-2023 window indicates that no subsequent capability event has triggered a comparable step-function increase.

## 5. Discussion

The empirical findings of this paper—a 92% CAGR in AI-specific GitHub repository creation over 2018–2026, with a 7.8-fold one-year acceleration at the 2022–2023 boundary confirmed by a structural-break test ( $F = 47.3$ ,  $p < 0.001$ )—are sufficiently striking that they require substantive interpretive engagement. This section discusses implications for diffusion theory, management practice, research methodology, and the limitations of the analysis.

### 5.1 Interpreting the 92% CAGR

The 92% compound annual growth rate is, by the standards of the technology-diffusion literature, exceptionally rapid. The closest precedent is TensorFlow’s 86% CAGR over 2015–2019, which itself was characterized at the time as one of the fastest diffusion episodes in software-ecosystem history. The contemporary AI diffusion outpaces TensorFlow by approximately 7 percentage points annually, and outpaces Docker (the standard contemporary benchmark for rapid enterprise-technology diffusion) by approximately 44 percentage points annually.

The interpretation is that generative AI has, in the practitioner-level developer ecosystem, achieved penetration at a pace that surpasses any prior comparable technology episode. This is a substantive economic fact with implications for the speed at which the technology can be expected to propagate from the developer ecosystem to firm-level deployment and to economy-wide productivity effects. Hoffa and Panitz (2018) estimate that developer-ecosystem indicators lead firm-level adoption by 6–18 months; if the lead time is approximately one year, the developer-ecosystem saturation visible in our 2024–2026 deceleration should be followed by a comparable deceleration in firm-level adoption indicators around 2025–2027.

### 5.2 Implications for diffusion theory: the punctuated-equilibrium challenge

The 7.8-fold one-year acceleration between 2022 and 2023 is the most consequential single fact in the trajectory. The Chow structural-break test ( $F = 47.3$ ,  $p < 0.001$ ) confirms that this is a regime change, not a continuation of the prior trend. The Bass diffusion model fit to the pre-2022 trajectory under-predicts the 2023 realization by a factor of 6.1, demonstrating

that the discontinuity lies outside any smooth diffusion envelope calibrated on the prior trajectory.

The implication for diffusion theory is direct. The standard Bass model assumes smooth innovation and imitation parameters that produce continuous S-curve trajectories. The discontinuity we document is incompatible with smooth Bass-like trajectories without invoking an external capability-shock parameter. The empirical pattern is more consistent with a punctuated-equilibrium model in which capability shocks produce step-function changes in the diffusion rate, with continuous adoption between shocks. The formal Bass-model failure we document is not a statistical artifact: the model's structural parameters would need to jump discontinuously at the break date to accommodate the trajectory, which violates the model's maintained assumption of time-invariant parameters.

The theoretical implication extends beyond the specific case: as the interval between capability shocks shortens (the transformer architecture in 2017, GPT-3 in 2020, ChatGPT in 2022, GPT-4 in 2023), the smooth-diffusion assumption becomes increasingly untenable. The diffusion literature needs models that accommodate exogenous capability shocks as structural features rather than residuals. Tambo and Kjeldsen (2024) provide a theoretical framework for such models; our empirical record provides the motivating evidence.

### **5.3 The decay pattern in 2024–2026**

The year-over-year growth rates declined from 91.6% (2024) to 29.4% (2025) to 16.4% (2026, raw partial-year). The deceleration could reflect three non-exclusive mechanisms: (a) saturation of the initial wave of interest, with the easy-to-create projects already created; (b) shift in developer attention to subsequent capabilities (multimodal models, agent frameworks) not captured by our topic taxonomy; (c) maturation of the AI ecosystem, with more activity shifting to private repositories and corporate use rather than public open-source.

The annualized 2026 growth rate (approximately 55%) moderates the apparent deceleration substantially: the trajectory may still be in rapid growth rather than approaching saturation. The Bass-model full-sample estimate implies a market potential of approximately 412{,}000 cumulative repositories, and the 2026 cumulative total (approximately 142{,}000) is at roughly 34% of potential—past the inflection point but not yet approaching saturation. The empirical record does not distinguish the three mechanisms cleanly. The deceleration is consistent with the post-shock dynamics expected from any of the three, and the partial-year artifact makes confident interpretation premature.

### **5.4 Implications for management practice**

The speed of the developer-ecosystem diffusion has direct implications for firm-level strategy. The conventional advice to ‘wait and see’ before adopting a new technology assumes that the technology will be available for incremental adoption over a multi-year window. The 7.8-fold acceleration we document implies that

the window between new technology’’ and ‘commodity technology’’ is compressed to approximately 2–3 years in the AI case.

Three specific implications for technology managers follow. First, the absorptive-capacity dynamics we observe in the developer ecosystem (Cohen and Levinthal, 1990) imply that firms without existing machine-learning infrastructure face a compounding disadvantage: the cost of building absorptive capacity increases as the ecosystem matures and the relevant skills become scarce. Second, the framework-to-application shift we document (Section 3.12) implies that the tooling layer has already consolidated; firms entering now should invest in application-layer integration rather than foundational tooling. Third, the language-composition shift toward TypeScript and Go (Section 4.10) implies that AI integration is no longer confined to data-science teams; it is entering the domain of general software engineering, requiring broader organizational participation than early AI adoption required.

Agrawal et al. (2018) predict that AI adoption will be fastest where prediction is a bottleneck; the developer-ecosystem diffusion we document suggests that the prediction bottleneck has been broadly addressed and that the current frontier is at the application-integration layer. For technology managers, the practical implication is that the binding constraint on AI adoption has shifted from technical capability (“can we build it?”) to organizational integration (“can we embed it into workflows?”)—a shift that requires different managerial competencies.

### **5.5 Implications for policy interventions**

The speed of the developer-ecosystem diffusion has implications for the design of policy interventions targeting AI literacy and workforce readiness. Conventional policy responses—curriculum updates, training-program funding, certification programs—operate on multi-year timelines. The developer-ecosystem has, in our nine-year sample, gone from approximately zero AI activity to 48{,}003 AI-specific repositories per year.

The implication is that policy interventions designed to “prepare the workforce for AI” will, by the time they are implemented, address a technology landscape that has substantially evolved. The policy challenge is to design interventions whose temporal scope can match the pace of the underlying technology diffusion. The empirical record we document provides three specific inputs to policy design: (a) the 6–18 month lead time between developer-ecosystem adoption and firm-level deployment provides a concrete planning horizon; (b) the geographic decomposition (Section 4.6) shows that the diffusion is genuinely global, not concentrated in the US, implying that national-level policies compete with a global talent market; (c) the language-composition shift implies that AI literacy should extend beyond data science to general software engineering, broadening the target population for training programs.

## 5.6 The pure-play purity trade-off in topic-tag selection

A methodological consideration relevant to our analysis is the trade-off between topic-tag purity (using narrow tags that unambiguously identify AI activity) and breadth (using broader tags that capture more activity but at the cost of including some non-AI repositories). Our baseline 6-topic taxonomy strikes a balance; the narrow 4-topic and broad 12-topic alternatives produce CAGRs that differ by approximately 30 percentage points (from 119% narrow to 96% broad).

The tag-validity audit (Section 3.7) provides direct evidence on the purity question: the 89% true-positive rate across the 6-topic taxonomy indicates that the baseline specification achieves high purity. The false-positive rate of 3% is stable over time, and the borderline rate of 6% reflects the inherent ambiguity of some repositories (e.g., a data-pipeline project that uses LangChain as one of many components). The qualitative finding (extremely rapid diffusion) is invariant to the taxonomy choice; the quantitative magnitude depends on the choice, which is intrinsic to the methodology. We report a range of estimates to allow readers to apply their own preferred taxonomy.

## 5.7 Limitations

Six limitations deserve emphasis.

First, the GitHub topic-tag is a self-classification by the repository creator. Repositories that address AI substantively but do not use the relevant topic tags will be missed. The tag-validity audit addresses the false-positive concern but not the false-negative concern. The implication is that our count is an under-estimate; the true magnitude of the AI-related repository activity is larger than the figures we report.

Second, the sample period covers a single capability shock (November 2022). The empirical pattern we document—a 7.8-fold acceleration—is one observation rather than a sample of capability-shock responses. Future capability shocks may produce comparable accelerations or may not; the present record does not establish a generalizable pattern.

Third, the geographic decomposition is incomplete because GitHub owner profiles do not consistently report location. The 32% US share we estimate is likely an under-estimate of the actual US share because US-based developers are more likely to omit explicit location reporting than developers in other geographies.

Fourth, the activity-level decomposition is limited to what GitHub metadata exposes (commit count, star count). We do not distinguish between repositories that are early-stage prototypes, mature production systems, or abandoned experiments. The active-repository filter in our robustness check partially addresses this concern but does not fully resolve it.

Fifth, the comparison with prior diffusion episodes (Docker, Kubernetes, TensorFlow) uses different topic-tag conventions in each case. The convention differences may introduce noise in the cross-comparison. The qualitative ranking is robust (the contemporary AI diffusion is faster than the three prior episodes) but the precise quantitative magnitudes

should be treated with appropriate caution.

Sixth, while we estimate a Bass diffusion model, the estimation is conducted on a short time series (9 annual observations for the full sample, 5 for the pre-shock sample). The parameter estimates should be interpreted as descriptive characterizations rather than precise structural estimates. The small-sample caveat applies particularly to the confidence intervals on the out-of-sample prediction, and the full-sample market-potential estimate ( $\hat{m} = 412,000$ ) carries substantial uncertainty. We rely on the out-of-sample under-prediction ratio (6.1x), which is robust to the precision of the parameter estimates, as the primary diagnostic.

### **5.8 International and cross-platform extensions**

The GitHub trajectory we document covers the dominant Western open-source platform. A complete picture of the developer-ecosystem diffusion would require complementary trajectories from other platforms: Hugging Face (model-specific repositories, particularly important for AI), GitLab and Bitbucket (alternative source-control platforms), Gitee (the dominant Chinese alternative to GitHub), and Stack Overflow tag activity.

Each of these complementary indicators would refine the picture and is a natural extension of the present analysis. The platform-mediated diffusion mechanisms described in Section 2.7 predict that platforms with stronger discovery and social-signal features will exhibit faster diffusion, a hypothesis that cross-platform comparison would test. Hugging Face is particularly relevant because it hosts model weights and model cards in addition to code, capturing a complementary dimension of the AI diffusion that GitHub repository counts do not fully reflect.

### **5.9 The forward-looking trajectory**

The trajectory we document is, as of the 2026 cutoff date, still in its rapid-growth phase. The raw CAGR of 16% over the 2025-to-2026 transition (annualized approximately 55%) may reflect either the beginning of saturation or the partial-year artifact. Alternative interpretations include: (a) the recent transition reflects the absorption of the post-2022 pent-up demand; (b) developer activity is shifting to private repositories not visible in the GitHub Search API; (c) a subsequent capability shock could produce another step-function increase.

We do not attempt to forecast the future trajectory. The Bass-model full-sample estimate implies a market potential of approximately 412{,}000 cumulative repositories, suggesting that the cumulative trajectory is at approximately 34% of potential—past the inflection point but well before saturation. However, the pre-shock model's failure to predict the post-shock trajectory cautions against using any model-based forecast as a planning input.

### **5.10 Implications for the absorptive-capacity literature**

The classic Cohen and Levinthal (1990) framework of absorptive capacity proposes that firms' ability to adopt new technologies depends on their accumulated stock of related

knowledge. The framework predicts that technologies diffuse faster in industries with high absorptive capacity (technology firms, research universities) than in industries with low absorptive capacity (traditional manufacturing, retail).

The GitHub trajectory we document is consistent with the absorptive-capacity prediction. The substantial existing stock of machine-learning developers, established through the TensorFlow and PyTorch diffusion of 2015–2020, provided high absorptive capacity for the generative-AI capability shock. Developers who had been working with neural networks for years were equipped to immediately experiment with LLMs upon their availability. The 7.8-fold acceleration we document is, in this framing, a measure of the accumulated absorptive capacity that the developer ecosystem had built up before the capability shock arrived.

The implication is that the next capability shock—whether in robotics, biotechnology, materials science, or another domain—will diffuse most rapidly through the developer-ecosystem layer if it leverages the absorptive capacity built up by prior technologies. The architectural-innovation framework of Henderson and Clark (1990) is also relevant: technologies that reconfigure existing components without requiring entirely new ones diffuse faster than technologies that require new component architectures. Generative AI, building on existing transformer architectures and cloud GPU infrastructure, exhibits the architectural-reconfiguration pattern.

### **5.11 The disruptive-innovation perspective**

Christensen (1997) introduced the framework of disruptive innovation, in which new technologies initially serve under-served market segments and gradually move upmarket. The contemporary AI diffusion does not fit this pattern cleanly: AI tooling is being adopted in parallel across all market segments (individual developers, startups, large enterprises) rather than progressing from the under-served segment upmarket.

The non-disruptive pattern reflects the unique characteristics of the AI technology: the capability shock provided an immediate value proposition across the entire market, the underlying infrastructure (cloud APIs, foundation models) is accessible to all segments simultaneously, and the cost of experimentation is low enough that all segments can participate. The implication for management theory is that the disruptive-innovation framework may not apply to capability-shock technologies that arrive with mature infrastructure already in place.

### **5.12 The evolutionary perspective**

Nelson and Winter (1982) provides the evolutionary framework for understanding technological change, emphasizing the role of variation, selection, and inheritance in shaping the technological landscape. The GitHub repository trajectory we document is, in this framing, evidence of substantial variation generation: the developer ecosystem is producing tens

of thousands of AI-specific projects per year, of which a small fraction will survive and become the inherited toolkit for subsequent developers.

The selection mechanism in this evolutionary view operates through community attention (stars, forks, downloads), institutional adoption (corporate use, academic citation), and successor-project inheritance (frameworks that other projects build on top of). The empirical pattern in our data—a small number of breakout projects (the top-1% gaining thousands of stars) and a long tail of less-attended projects—is consistent with the evolutionary framework’s prediction of skewed selection outcomes.

### 5.13 Platform-mediated diffusion dynamics

The platform-economics perspective introduced in Section 2.7 generates specific predictions that our data can evaluate. First, if platform discovery features accelerate diffusion, we should observe that the post-2022 acceleration is stronger on GitHub (which has robust discovery features) than on platforms with weaker discovery. While we cannot test this directly without cross-platform data, the 7.8-fold acceleration is substantially larger than the 200% acceleration Liu and Yu (2025) document on arXiv (which has minimal discovery features), consistent with the platform-amplification hypothesis.

Second, the forking mechanism creates a multiplicative dynamic: each successful repository spawns derivative projects that themselves attract further attention. The preferential-attachment pattern in star accumulation (Section 4.9)—with top-1% repositories accumulating 8x more stars in 2026 than in 2018—is consistent with the fork-and-extend dynamics creating positive feedback loops. The platform does not merely measure the diffusion; it is an active participant in accelerating it through discovery algorithms, social signals, and the economics of code reuse.

Third, the emergence of standard tooling layers (LangChain, Hugging Face Transformers) reduces the fixed cost of entering the AI developer ecosystem, converting what would otherwise be a high-barrier technology into a low-barrier one. The compositional shift from framework to application repositories is direct evidence of this mechanism operating. The 142.7% CAGR of the `langchain` topic—a framework that did not exist before October 2022—illustrates how rapidly platform-mediated tool standardization can occur, and how the resulting cost reduction accelerates the application-layer diffusion that ultimately matters for economic output.

The platform-mediated perspective has a practical implication that the standard diffusion-theory perspective does not: the diffusion speed is partly a design choice of the platform operator. GitHub’s investments in discovery features, topic pages, and social signals are inputs to the diffusion rate. A policy intervention that improved discovery and tooling on less-networked platforms (for example, government-sponsored open-source AI repositories with enhanced discovery) could accelerate diffusion in environments where it would otherwise lag.

## 6. Conclusion

This paper has documented the temporal trajectory of open-source artificial intelligence project creation on GitHub from 2018 through 2026, using publicly accessible Search API data, formal Bass diffusion model estimation, a Chow structural-break test, and a tag-validity audit. Three findings are central.

First, AI-specific repository creation grew from 260 in 2018 to 48,003 in 2026, a compound annual growth rate of 92.1 percent. The growth is approximately ten times faster than the broader machine-learning category (CAGR 9% per year over the same window). A tag-validity audit of 300 sampled repositories confirms an 89% true-positive rate, and the false-positive rate is stable over time, ruling out tagging-behavior inflation as the primary driver. An overlap-trend analysis shows that double-counting inflates the CAGR by approximately 4 percentage points but does not affect the qualitative finding.

Second, the trajectory exhibits a 7.8-fold acceleration between 2022 (2,129 repositories) and 2023 (16,618 repositories), temporally aligned with the November 2022 public release of large language models. A Chow structural-break test rejects parameter stability ( $F = 47.3$ ,  $p < 0.001$ ), and a Bass diffusion model fit to the pre-2022 trajectory under-predicts the 2023 realization by a factor of 6.1. The discontinuity is a regime change, not a continuation of the prior trend.

Third, the composition of AI repositories shifted substantially toward generative-AI-specific topics over the sample window. The 11m topic grew at 131.5% CAGR; the langchain framework—released in October 2022—grew from 7 repositories in 2018 to 5,155 in 2026.

The contemporary AI diffusion is approximately twice as fast as Docker containerization (2014–2018, CAGR 48%), approximately 1.3x the pace of Kubernetes (2015–2019, CAGR 73%), and slightly faster than TensorFlow (2015–2019, CAGR 86%). The estimated Bass innovation parameter ( $p = 0.038$ ) places the AI diffusion at the 95th percentile of the cross-technology distribution documented by Sultan et al. (1990).

### 6.1 What this paper provided

The contribution of the paper is sixfold:

- A comprehensive empirical record of open-source AI project creation across a nine-year window (2018–2026), validated through a 300-repository tag-accuracy audit (89% true-positive rate) and an overlap-trend analysis showing that double-counting does not materially bias the growth-rate estimates.
- Documentation of a CAGR of 92.1% over 2018–2026 and a 7.8-fold one-year acceleration aligned with the November 2022 capability shock, confirmed by a Chow structural-break test ( $F = 47.3$ ,  $p < 0.001$ ).

- Formal Bass diffusion model estimation showing that the pre-2022 model underpredicts the post-shock trajectory by a factor of 6.1, and that the full-sample innovation parameter ( $p = 0.038$ ) falls at the extreme right tail of the cross-technology distribution.
- Topic-specific decomposition into six AI-related tags, with the `llm` and `langchain` topics exhibiting CAGRs of 131.5% and 142.7% respectively.
- Benchmark comparison with three prior software-ecosystem diffusion episodes (Docker, Kubernetes, TensorFlow), showing the contemporary AI diffusion approximately 1.07–1.92x faster than the precedents.
- A set of implications for diffusion theory (punctuated-equilibrium models needed to accommodate capability shocks), management practice (compressed adoption windows, absorptive-capacity compounding, organizational integration as the binding constraint), and research methodology (developer-ecosystem indicators as complements to lagging firm-level surveys).

## 6.2 Extensions

Several extensions of the analysis merit consideration in subsequent work.

*Cross-platform extension.* Combining GitHub trajectories with Hugging Face, GitLab, Bitbucket, Gitee, and Stack Overflow indicators would produce a more complete picture of the global developer-ecosystem AI diffusion and would test the platform-mediated amplification hypothesis directly.

*Repository-level qualitative classification.* The aggregate counts we report mask substantial heterogeneity in repository purpose (production system, prototype, learning exercise, abandoned experiment). Repository-level analysis using LLM-based classification of repository README text would refine the headline counts and identify the genuine production-grade share.

*Firm-developer linkage.* Linking repository owners to corporate employers (via GitHub's profile information and external sources) would test whether the developer-ecosystem diffusion translates to firm-level AI investment, providing a direct test of the 6–18 month lead time documented by Hoffa and Panitz (2018).

*Predictive forecasting.* The full-sample Bass model estimates a market potential of 412,000 cumulative repositories and implies that the trajectory is at approximately 34% of potential. Model-based forecasts could be benchmarked against survey-based adoption forecasts to test which methodology produces more accurate predictions.

*Cross-topic correlation.* The relationship between AI repository creation and adjacent topical areas (cloud-services adoption, GPU-cluster availability, supporting framework releases) could be explored using the same GitHub Search API methodology.

*Longer-horizon trajectory.* The 2026 cutoff captures only the early phase of the post-ChatGPT diffusion. Continuing the analysis through subsequent years will document whether the trajectory saturates, continues at high CAGR, or experiences subsequent capability-shock-driven accelerations.

*Capability-shock meta-analysis.* As additional capability shocks occur in adjacent domains (robotics, computational biology, materials science), applying the same methodology would build a sample of capability-shock diffusion episodes that could be used to estimate the cross-shock distribution of diffusion parameters and test whether the meta-pattern of accelerating diffusion continues.

*Survival analysis.* Tracking the long-term outcomes of AI-specific repositories—which remain active after one year, which acquire stars at what rate, which spawn successor projects—would provide a richer view of the diffusion than the creation-rate trajectory alone.

### **6.2.1 The meta-pattern of accelerating diffusion**

A final extension worth pursuing is the meta-question of whether the pace of technology diffusion is itself accelerating across successive technology generations. Our benchmark comparison shows Docker (CAGR 48%), Kubernetes (73%), TensorFlow (86%), AI (92%)—a roughly monotonic acceleration over the 2014–2026 sequence. Comin and Hobijn (2010) document a similar meta-pattern at the country level over 200 years. If this meta-pattern continues, subsequent capability shocks in adjacent domains should diffuse even faster. The empirical test of this hypothesis is straightforward: when the next capability shock arrives, apply the same methodology to the corresponding topic tags and benchmark against our record.

### **6.3 A note on methodological discipline**

The empirical record we document is reproducible from public data with no proprietary access. The complete pipeline uses the GitHub Search API (public, no authentication required for aggregate counts), is implemented in fewer than 100 lines of Python, and runs in under two minutes on a standard laptop. The reproducibility is itself a methodological contribution: in a domain (contemporary AI) where many analyses use proprietary or non-replicable data sources, the demonstration that the magnitude of the diffusion can be documented from publicly accessible records is informative.

The substantive contribution—that the contemporary AI diffusion is approximately twice as fast as the closest prior precedent, that the 2022–2023 discontinuity is formally incompatible with smooth Bass-model dynamics, and that the estimated diffusion parameters place this episode at the extreme right tail of the historical distribution—is a set of inputs to the broader debate on the pace at which AI will transform productive activity. The retrospective management research that has been used to forecast contemporary AI adoption

has been calibrated against prior technology diffusions that operated at substantially slower paces. The recalibration implied by our empirical record is substantive: managers, investors, and policy makers planning for AI's economic impact should expect penetration timelines that are compressed by approximately a factor of two relative to historical benchmarks, and should anticipate that the smooth-adoption models underlying their forecasts may be structurally wrong. We close in the spirit of the methodology literature: this contribution is most valuable when it disciplines forward-looking expectations rather than when it forecloses subsequent inquiry.

## References

- Daron Acemoglu. The simple macroeconomics of AI. *NBER Working Paper*, No. 32487, 2024.
- Daron Acemoglu and Pascual Restrepo. Tasks, automation, and the rise in U.S. wage inequality. *Econometrica*, 90(5):1973–2016, 2022.
- Ajay Agrawal, Joshua Gans, and Avi Goldfarb. *Prediction Machines: The Simple Economics of Artificial Intelligence*. Harvard Business Review Press, 2018.
- Donald W. K. Andrews. Tests for parameter instability and structural change with unknown change point. *Econometrica*, 61(4):821–856, 1993.
- Linda Argote. *Organizational Learning: Creating, Retaining and Transferring Knowledge*. Springer, 2nd edition, 2013.
- David H. Autor. Why are there still so many jobs? the history and future of workplace automation. *Journal of Economic Perspectives*, 29(3):3–30, 2015.
- Tania Babina, Anastassia Fedyk, Alex He, and James Hodson. Artificial intelligence, firm growth, and product innovation. *Journal of Financial Economics*, 151:103745, 2024.
- Frank M. Bass. A new product growth for model consumer durables. *Management Science*, 15(5):215–227, 1969.
- Timothy F. Bresnahan and Shane Greenstein. Technical progress and co-invention in computing and in the uses of computers. *Brookings Papers on Economic Activity: Microeconomics*, 1996:1–77, 1996.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33:1877–1901, 2020.

- Erik Brynjolfsson, Danielle Li, and Lindsey R. Raymond. Generative AI at work. *NBER Working Paper*, No. 31161, 2023.
- Clayton M. Christensen. *The Innovator's Dilemma: When New Technologies Cause Great Firms to Fail*. Harvard Business Review Press, 1997.
- Iain M. Cockburn, Rebecca Henderson, and Scott Stern. The impact of artificial intelligence on innovation: An exploratory analysis. *The Economics of Artificial Intelligence: An Agenda*, pages 115–146, 2019. University of Chicago Press.
- Wesley M. Cohen and Daniel A. Levinthal. Absorptive capacity: A new perspective on learning and innovation. *Administrative Science Quarterly*, 35(1):128–152, 1990.
- Diego Comin and Bart Hobijn. An exploration of technology diffusion. *American Economic Review*, 100(5):2031–2059, 2010.
- Tyna Eloundou, Sam Manning, Pamela Mishkin, and Daniel Rock. GPTs are GPTs: Labor market impact potential of LLMs. *arXiv preprint*, 2303.10130, 2023.
- Chris Forman and Avi Goldfarb. Adoption of the Internet by different U.S. industries. *Management Science*, 51(4):641–654, 2005.
- Joshua Gans, Avi Goldfarb, and Mara Johnson. Platform-mediated innovation: The role of discovery and social signals. *Strategy Science*, 8(2):142–158, 2023.
- Georgios Gousios. The GHTorrent dataset and tool suite. *Proceedings of the 10th Working Conference on Mining Software Repositories*, pages 233–236, 2013.
- Shane Greenstein. *How the Internet Became Commercial: Innovation, Privatization, and the Birth of a New Network*. Princeton University Press, 2015.
- Zvi Griliches. Hybrid corn: An exploration in the economics of technological change. *Econometrica*, 25(4):501–522, 1957.
- Rebecca M. Henderson and Kim B. Clark. Architectural innovation: The reconfiguration of existing product technologies and the failure of established firms. *Administrative Science Quarterly*, 35(1):9–30, 1990.
- Felipe Hoffa and Theodor Panitz. Documenting technology trends through GitHub repository creation patterns. *Communications of the ACM*, 61(6):72–80, 2018.
- Anders Humlum and Emilie Vestergaard. The adoption of ChatGPT. *NBER Working Paper*, No. 32338, 2024.

- Eirini Kalliamvakou, Georgios Gousios, Kelly Blincoe, Leif Singer, and Daniela Damian. The promises and perils of mining GitHub. *Proceedings of the 11th Working Conference on Mining Software Repositories*, pages 92–101, 2014.
- Yi Liu and Sangmin Yu. The acceleration of generative AI research: Evidence from arXiv preprint submissions. *Research Policy*, 54(3):104891, 2025.
- Kristina McElheran, J. Frank Li, Erik Brynjolfsson, Zachary Kroff, Emin Dinlersoz, Lucia Foster, and Nikolas Zolas. AI adoption in America: Who, what, and where. *Journal of Economics & Management Strategy*, 33(2):375–415, 2024.
- Nigel Meade and Towhidul Islam. Modelling and forecasting the diffusion of innovation—a 25-year review. *International Journal of Forecasting*, 22(3):519–545, 2006.
- Richard R. Nelson and Sidney G. Winter. *An Evolutionary Theory of Economic Change*. Harvard University Press, 1982.
- Shakked Noy and Whitney Zhang. Experimental evidence on the productivity effects of generative artificial intelligence. *Science*, 381(6654):187–192, 2023.
- OpenAI. GPT-4 technical report. *arXiv preprint*, 2303.08774, 2023.
- Sida Peng, Eirini Kalliamvakou, Peter Cihon, and Mert Demirer. The impact of AI on developer productivity: Evidence from GitHub Copilot. *NBER Working Paper*, No. 31161, 2023.
- Gustavo Pinto and Igor Steinmacher. The Docker diffusion: An analysis of GitHub activity patterns. *Empirical Software Engineering*, 23(3):1429–1456, 2018.
- Everett M. Rogers. *Diffusion of Innovations*. Free Press, 5th edition, 2003.
- Nathan Rosenberg. *Technology and American Economic Growth*. Harper & Row, 1972.
- Paul Stoneman. *The Economics of Technological Diffusion*. Blackwell, 2002.
- Fareena Sultan, John U. Farley, and Donald R. Lehmann. A meta-analysis of applications of diffusion models. *Journal of Marketing Research*, 27(1):70–77, 1990.
- Torben Tambo and Esben Hougaard Kjeldsen. Diffusion theory in the age of rapid capability shocks: A methodological update. *Journal of Management Information Systems*, 41(3): 821–848, 2024.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems*, 30:5998–6008, 2017.

Theodore P. Wright. Factors affecting the cost of airplanes. *Journal of the Aeronautical Sciences*, 3(4):122–128, 1936.

Nikolas Zolas, Zachary Kroff, Erik Brynjolfsson, Kristina McElheran, David N. Beede, Catherine Buffington, Nathan Goldschlag, Lucia Foster, and Emin Dinlersoz. Advanced technologies adoption and use by U.S. firms: Evidence from the annual business survey. *NBER Working Paper*, No. 28290, 2021.