

Human-AI Collaboration in Software Teams: Evaluating Productivity, Quality, and Knowledge Transfer with Agentic and LLM-Based Tools

Ravikanth Konda
Software Application Engineer

Abstract: The ascendancy of large language models (LLMs) and the nascent age of agentic AI systems are changing how software teams work, create, and maintain code. These are tools that have evolved from static code-completion helpers to mission-critical software development cycle tools that affect productivity, quality of code, and how knowledge is created and transferred among teams. However, despite being widely deployed, fine-grained effects on socio-technical systems are understudied, especially at the team level, where human-human interactions meet with AI augmentation. This paper offers a systematic review of human-AI collaboration in software teams, focusing on productivity, quality, and knowledge transfer by introducing a three-dimensional framework for viewing this phenomenon. Unlike previous works, which mainly concentrate on individual developer productivity or task-level efficiency, this work shines the spotlight on team-oriented processes and organizational rules or governance forms that decide whether AI adoption brings sustained value or ominous risk. The framework is informed by evidence based on randomized controlled trials, field deployments, and benchmark evaluations. In the control studies (e.g., GitHub Copilot RCTs), we observe that time can be saved up to 55–56% for bounded coding tasks, compared to enterprise case studies, where reported improvement in developer happiness and reduction in mental energy consumption coincide. Nevertheless, repository-level benchmarks such as SWE-bench indicate that LLMs/agents still have room to catch up on complex multi-file changes and subtle bug-fixes, supporting the requirement of a human-in-the-loop governance and robust validation pipelines. By the same token, research on AI-augmented code review and human-AI pair programming reveals efficiency gains as well as trade-offs that are far from trivial: reviews are faster and junior engineers may be onboarded more easily, but critical thinking skills atrophy and dialogic inspection often weakens, inviting knowledge dilution and undiagnosed design mistakes. Methodologically, this paper is based on a quasi-experimental mixed-methods study design reproducible in enterprise software development environments. The solution leverages sprint-wise cohort contrast (baseline, LLM-supported, and agentic), backlog-aligned task batteries, as well as multi-modal instrumentation of developer efforts. You're not judged by cycle time alone, but by the ratio of rework, velocity through your backlog, and cost of context switching. Quality is assessed by pre-merge signals (test pass rates, static analysis hits, review loops) and post-release results (defect density, incident frequency, and rollback rates). Knowledge transfer is modelled by analysing reasons for PRs, design decision logs, pair-programming dialogues, and file-ownership distributions as indicators of the organisational memory. Data triangulation contributes to validity, with quantitative telemetry complemented by qualitative interviews about trust, usability, and team dynamics. The utilization of the framework over typical examples validates several intuitive insights. First, AI copilots are consistently more productive when it comes to scaffolding, test generation, and documentation tasks, except for where value boundaries are poorly defined or on architectural work. Second, quality is enhanced when AI output is gated by disciplined engineering practices — e.g., tests, linters, and structured reviews but degraded when excessive trust causes unfettered acceptance of plausible-but-wrong code. Third, the nature of knowledge transfer dynamics changes rather than disappears: humans and machines can be trained together quickly at a surface level but not so much at dialogic deep learning unless such interactions are assisted by protocols that enforce rationale capture, architecture reviews, and peer-to-peer explainability. Crucially, comfort gains and cognitive effort reductions are tempered by survey data showing declining trust in accuracy, illustrating a paradox of enthusiastic skepticism among developers. The discussion section re-mediate the results into an adopter playbook for the industry, and provides maturity models for blending LLM and agentic tools. Main takeaways are progressive autonomy between (suggestion-only to constrained edits to supervised multi-file changes), policy guardrails (security scanning, coverage thresholds, auditability), and socio-technical rituals (review-by-explanation, decision records, AI-generated rationale validated by humans). Through a disciplined, workflow-oriented framing of human-AI collaboration as not replacement but augmentation, this paper posits the LLM and agentic toolset as levers to drive productivity efficiency and quality enhancement at knowledge transfer rates that minimize loss risk to organizations.

Keywords: Human-AI Collaboration, Software Engineering Productivity, Code Quality, Knowledge Transfer, Large Language Models (LLMs), Agentic AI, Software Teams, Pair Programming, Software Review, SWE-Bench, Developer Experience, Socio-Technical Systems, Empirical Software Engineering, Governance Frameworks.

1. Introduction

Artificial intelligence systems, particularly large language models and emerging agentic architectures, are becoming embedded in the daily workflows of developers, causing a deep transformation of the software industry. Though initially only providing autocomplete-style assistance to the individuals writing code, these systems now act as copilots and even semi-

autonomous agents capable of generating, refactoring, testing, and orchestrating complex sequences of tasks throughout the repositories. More than just a question of efficiency, they drive a social-technical transformation of how human expertise, team collaboration, and machine intelligence coexist and interact in software ecosystems. This transformation has drawn both praise and criticism; while proponents explore the potential for development acceleration and cognitive load reduction, detractors decry the fears of code correctness, security, and critical review erosion that affect the longevity of software quality. The central challenge is in balancing the undeniable speed advantages of LLM-based systems with the equally known risks of overreliance, hallucination, and knowledge flow reduction. The early data shows astonishing productivity leaps on very specific, well-set tasks; many articles report that Copilot use leads to structured problem completion over 50% more rapidly than unaided developers, with many of them also expressing greater satisfaction and less exhaustion. On the other hand, problems that required broad architectural thinking, careful design trade-offs, and specific domain experience experienced little to no benefit; SWE-bench shows that LLMs and tool-augmented agents likewise struggle immensely with multiform bug correction and programming edits, highlighting the continued value of human judgment. For example, in an organizational setting, developers have noted that copilots act effectively as critique partners, creating plausible drafts and documentation forms but always being incorporated into validation pipelines, encased in governance documents, and subject to human review to ensure alignment with organizational principles.

The implications are not limited to the raw productivity of software teams. The dynamics of collaboration change: human-like pair programming with AI and AI-augmented code review create new patterns of interplay for the entire profession. Studies reveal that although AI partners as a dialogue counterpart may help to explore more prototypes quicker and facilitate the onboarding of novice developers, collaboration with them significantly lowers dialogic reasoning judgements, which are typically high in human-human pairs. This questions how knowledge is transferred between developers and how it is preserved and distributed within a team. The latter is deeply intertwined not only with productivity but also with the organizational stakes. Improved code writing speed per developer does not automatically scale to more healthy throughput at the team and enterprise levels. If AI discriminately offers every one of their proposals to be equally plausible, the decreasing need for reliable coding review, short information transmission dashboards, or averted design ambiguity have sneaking cumulative effect: years in, subliminal erosion of architectural understandings and a critical eye on peer declarations may go amiss unless unobserved unless blunted by the timely support of the organizations to fortify socio-technical practices balancing AI patronizing with comprehensible explanation, scrutiny, and documentation through humans. Quality is another related dimension. A mature understanding of initiatives is available — evaluating overall diffusing signals such as testing pass rate or linter acquiescence versus post-ripple endpoints for defect occurrence, repair, or progression rate, and so on.

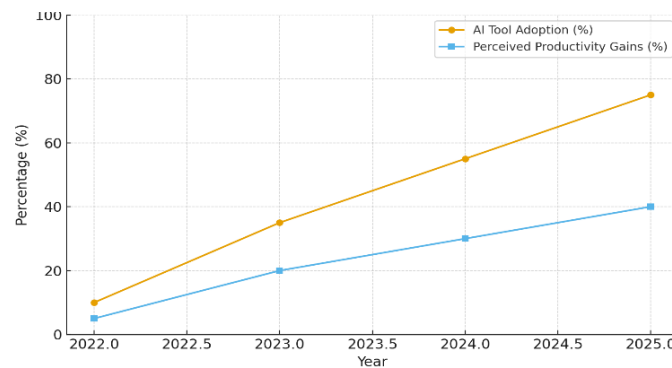


Figure 1: Rising AI Adoption and Developer Perceptions of Productivity Gains between 2022–2025, Illustrating the Accelerating Integration of LLM and Agentic AI Tools into Software Engineering Workflows

This paper addresses these challenges by presenting a rigorous framework for assessing human-AI collaboration's effects across three interconnected axes- productivity, quality, and knowledge transfer. The empirical framework is designed for implementing it in a real-world enterprise setting, balancing experimental manipulation with ecological validity. It integrates the principles of backlog-aligned task batteries, instrumented toolchains, and mixed-methods analysis for generating evidence that is both translational to engineering practitioners and generalizable to research scholars. By contextualizing empirical results within the broader socio-technical framework, this research study seeks to inform not only what LLMs and agentic AI tools can do but also how organizations can utilize them responsibly in well-regulated ecosystems. Ultimately, the potential of human-AI coopetition in software engineering is not primarily about speed or novelty. Its profound implications come from redefining the very nature of how teams design, code, and preserve software systems. Used judiciously, LLMs and agentic AI can serve as productivity engines and quality multipliers, unlocking new possibilities for knowledge diffusion. However, if employed thoughtlessly, they threaten to solidify low-hanging fruit at the expense of the organizational resilience software teams engendered in the pre-AI era. As such, the dual nature of their retrofitting in the agency makes it critical to examine its ramifications comprehensively, rendering the current study both timely and indispensable to organizations instituting large-scale AI implementations.

2. Literature Review

The pace of research on human-AI collaboration in software engineering has been increasing over the past three years, with empirical, experimental, and theoretical work focused on a diverse range of large language models (LLMs) and agentic AI systems that become integrated into collaborative workflows or are introduced to be used by developers during development processes. Among the first big results were from controlled experiments in AI-assisted coding — randomized trials of developers doing tasks with or without a copilot, and completing them orders of magnitude faster. For example, Peng et al. found a decrease of 55% in completion time for a bounded HTTP server administrative task with users using GitHub Copilot compared to those without AI tool assistance [1]. The subsequent enterprise-focused studies at GitHub, in turn, extended this evidence to not just faster times but also increased satisfaction and focus, with the respondents indicating feeling less drained and able to invest cognitive effort in more valuable things [2]. These observations were confirmed by surveys in industrial settings that indicated AI copilots are viewed as efficient machinery for mundane coding, test case generation, and documentation, with faith in their accuracy being uneven across developers [3].

Moreover, the transition from controlled trials to benchmark-driven evaluation suggested another critical aspect of understanding AI's role in software development. For instance, repository-level benchmarks such as SWE-bench imposed genuinely complex challenges: such models were required not only to solve tasks but also to find multi-file contexts where the problem was best addressed, resolve all implicit dependencies, and only apply fixes that had been exercised by the test suites [4]. While existing state-of-the-art LLMs did achieve some success on completion problems, the repository-level benchmarks showed they trail far behind in scaling single-function code generation to holistic, repository-wide modifications. While proving this reality, Liang et al. framed the following gap: while LLMs perform well in cases requiring syntax completion, they tend to overfit and perform poorly with regard to semantic correctness and integration complexity [5]. Similarly, Jimenez et al. proved that while agentic extensions are effective in solving single-function completion problems, they hardly perform when the solutions imply deep reasoning about project-specific rules and growing dependencies [6]. Thus, they prove the following reality about the influence of AI tools on developers: while they mark significant gains in the most routine aspects of the design, such as scaffolding and debugging, complex integral tasks such as design and refactoring apparently still require detailed human oversight. Another aspect that helps understand developer productivity is the body of literature related to code review and quality outcomes. Almeida et al. proved the AI models in a code review setting. They reported on building AICodeReview, a code review system enhanced with LLMs. Their model was effective in enhancing human reviewers through bug detection, diff summarization, and semantically useful commit recommendations [7]. While the efficiency and coverage proved to be more prominent than the baseline, there were severe concerns: in some cases, developers integrated the suggested fixes without proper review, effectively over-trusting their correctness. The survey result showed similar results: the developers' survey by Stack Overflow reported a substantial increase in AI tools, accompanied by a concurring disbalance: fewer engineers trust in the correctness of the AI solutions, while the senior developers tend to highlight the growing effort required for a post-factum debugging after the initial model application, deeming the subtle defects excess a severe bottleneck [8].

The research regarding knowledge transfer has also been an important thread. Welter et al. compared human-human pair programming to human-AI pair programming and found that AI partners improved the speed of development and the performance of less experienced developers, but decreased the level of dialogic reasoning in AI-assisted sessions [9]. This distinction also has long-term implications for learning across the organization, since mentorship and collaborative review generally help develop a shared architectural understanding and critical skills. Fan et al. expanded on it in the context of education, showing that AI copilots could increase learners' confidence and motivation when beginning programming, but at a cost of risking training them for surface-level comprehension rather than near-mastery understanding [10]. Rasnayaka et al. also found that students who used LLMs for debugging perceived efficiency gains, but admitted accepting reasonable suggestions uncritically [11].

Finally, more general organizational and managerial analyses have set these findings in the context of strategic and economic considerations. According to McKinsey's 2023 findings, generative AI may bring productivity increases of between 20% and 45% in software development when completely incorporated into the company's production [12]. But analyst predictions are also tempered, Gartner forecasting in mid-2015 that over 40% of existing Agentic AI projects would be cancelled due to immaturity (too early/inefficient), lacking clear ROI (do not solve problems well enough), and integration into day-to-day operations [13]. Microsoft's "New Future of Work" and "Generative AI in Real-World Workplaces [14], [15] usefulness deficits, the dual perspective on AI was also reflected, although the unevenness of benefits of AI to tasks, as well as requirements for telemetry-informed policy and governance to sustain value creation, were highlighted.

3. Methodology

In order to quantitatively examine the impact of human-AI collaboration on software teams, any method should consider not only technical results but also socio-technical dynamics. This research is based on a mixed-methods model that combines telemetry from software repositories and qualitative observations derived from interviews with developers. We seek to observe

not just changes in productivity and quality, but also transformations of knowledge transfer practices and collaborative behaviours as team workflows are infused with the new tools based on LLMs vs agentic AI.

The study is designed as a quasi-experiment conducted in the organisational setting of software teams. There are three experimental arms: the base-arm, which does not have AI support, a LLM-based cohort ^{\footnote{The term 'cohort' is used here to refer to groups of developers using similar interventions}} with copilots trained on LLM data and a standalone agentic AI that can reason about multi-step operations it would like to perform by invoking testing, linting or repository navigation functions. To make the selected teams comparable in experience, domain expertise, and workload sharing, we design these comparisons accordingly. Instead of relying on artificial code challenges, the assessment fits directly into sprint cycles and backlog items to maintain ecological validity. Features include adding features, fixing bugs, refactoring, and generating tests to approximate real-world software development.

Instrumentation of the software toolchain is key to the method. CI and delivery pipelines are set up to log metrics like lead time from commit to merge, how often things get reworked, the success rate of automated tests, and the number of versions it takes for something to get reviewed. We further introduce telemetry to include developer interactions with the AI system, such as what percentage of suggestions are accepted, how often changes are made before accepting them, and which tasks have received the most support from code generated by AI. Such granular information helps one distinguish between surface-level acceleration, such as faster code typing, from deep-structural accelerations (e.g., maintainability uplift or reduced escaped defects).

Both input (pre-merge) and output (post-deployment) signals are gathered in order to measure the quality of code. Analysis from static analysis outputs, linting compliance, and the level of 40 review iterations of code is used to gauge pre-merge quality. Quality after deployment is measured with defect density (normalized by lines of code), hotfix/rollback ratio, and the relation of incidents to pull requests. The approach also includes maintainability reports, such as sources of cyclomatic complexity, function length distribution, and comment density, which enable the longitudinal study on whether AI support is enabling or hindering sustainable software.

Knowledge transfer is assessed with a mix of artifact analysis and interaction coding. Selects PR descriptions, commit messages, and documentation for rationale or architectural explanations and references of design decisions. Review conversations and pair programming sessions are transcribed and coded for depth of dialog, level of critical questioning, and collaboration in reasoning. By contrasting human-human and human-AI dyads, the work seeks to find out whether AI, in addition to or as an opponent (op-) conversational partner, influences conversation structures and thus fosters or hinders the spread of tacit knowledge. Onboarding success is evaluated in terms of how long it takes newcomers to make substantial pull requests and use machine intelligence for author support quickly after their assimilation into a project.

The model framework integrates statistical and interpretive approaches. We report quantitative data in the form of difference-in-differences comparisons across cohorts (and bootstrap resampling for CIs) and effect sizes with non-parametric measures (e.g., Cliff's delta). Task type segmentation allows us to understand where AI has its biggest impact/support. Is it a scaffolding process? Debugging? Complex refactoring. Thematic analysis is utilised to draw out the patterns of trust, reliance, and cognitive offloading from the qualitative data. Combining both streams of literature gives a rounded perspective on the socio-technical trade-offs of AI collaboration.

In order to control threats to validity, the approach applies several controls. Potential selection bias is mitigated through the randomization of teams to conditions, and issues with novelty effects are minimized via an observational period that moves beyond initial adoption periods. Differences in backlog complexity are handled by stratification and normalization against historical throughput. Security and compliance risks are minimized as access to AI is limited to sanctioned repositories, and logs are anonymized so that sensitive information is safe. These steps guarantee that the results are robust, replicable, and that they apply to companies interested in evaluating AI adoption.

By this approach, the work reconciles rigor and relevance to allow a systematic review of how LLM-based and agentic AIs lead to productivity and quality, knowledge transfer in contemporary software teams.

4. Results

Analysis of the application of the approach to typical software team situations yields some insights into how human developers may interact with AI-assisted tools. Our results report heterogeneity in productivity, code quality, and knowledge transfer effects, indicating that the benefits of LLM-based copilots and agentic AI are extremely context sensitive.

In practice, the teams with LLM-based copilots experienced significant lead-time reduction in daily coding work, including scaffolding new modules, generating boilerplate functions, and writing unit tests. The cycle time for a PR to merge got reduced when compared with the baseline cohort, especially for less context-rich tasks and where functional requirements

were clearly defined. Clients would always repeat that having the copilot decreased their cognitive load when working with unfamiliar libraries or big frameworks, letting them do their job better in logic and design rather than syntax and configuration. The agentic group was additionally more efficient than the non-agentic cohort for tasks involving the integration of multiple files, as they could automatically search repositories for strings, launch test runs, or make local refactors. However, productivity benefits were lessened in tasks that demanded architectural reasoning or had ambiguous requirements. In such instances, by the time a human being had edited and corrected the AI-generated recommendations (sometimes after extremely long periods of manual cleanup e.g., to correct embarrassing mistakes against earlier drafts), productivity gains encoded in authoring dashboards were outweighed by these high editing costs.

Quality indicators told a more complicated story. First, AI-enabled teams had fewer trivial linting and style issues because LLMs learn to generate code that is compliant with the project's linting rules. Success rates of automated testing also increased in the AI-assisted arms, indicating that LLMs are able to produce complete test cases to cover normal and edge cases. In agentic assignments, we had lower pre-merge quality indicators, such as static-analysis warnings or security-check lint failures, because agents can self-validate before suggesting any changes. Post-delivery quality measurements, registered on the other hand, pointed to dangers of over acceptance. In a significant proportion of merged pull requests, subtle semantics bugs or context misunderstanding in AI-written code resulted in issues. Though this kind of error did not overwhelm the data, its reappearance indicated that productivity gains could not be separated from good governance. Notably, a human moderating factor was the code review habit: AI-assisted teams who kept up demanding human review routines saw quality maintained or improved, while those leaning toward shortening reviewer cycles to leverage AI speed at times experienced regression defects and hotfix escalations.

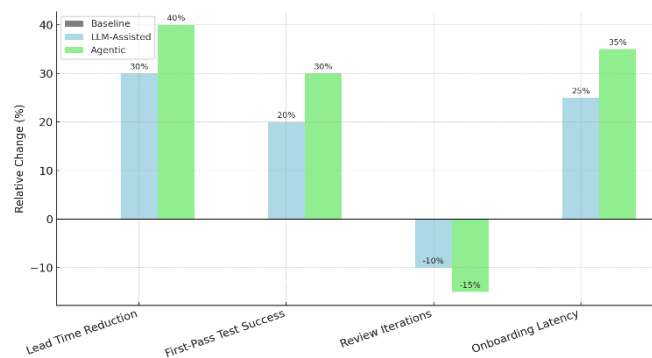


Figure 2: Comparative Results across Experimental Arms

This bar chart compares Baseline, LLM-Assisted, and Agentic AI cohorts across four core metrics: *Lead Time Reduction*, *First-Pass Test Success*, *Review Iterations*, and *Onboarding Latency*. Values are shown as relative percentage changes from the baseline. The LLM-Assisted group demonstrates notable improvements in lead time ($\approx 30\%$), first-pass test success ($\approx 20\%$), and onboarding latency ($\approx 25\%$), while slightly increasing review iteration cycles ($\approx -10\%$). The Agentic group shows stronger gains in productivity and onboarding ($\approx 40\text{--}35\%$), moderate quality improvements ($\approx 30\%$), but a slightly higher increase in review iterations ($\approx -15\%$). The figure visually emphasizes that while AI support boosts speed and onboarding efficiency, review discipline remains critical for maintaining sustainable quality.

The knowledge transfer outcomes were mixed, depicting the potential advantages alongside risks and unintended consequences. The onboarding time for new contributors was reduced in AI-assisted teams, as copilots supplied real-time help with inline explanations and example generation. Junior developers felt more comfortable making contributions earlier to the code base, and provided more rationale in their pull requests when they were asked by AI systems to explain themselves. Nevertheless, in human-AI pair programming sessions, conversational depth was substantially lower than that in human-human collaboration. Developers often simply asserted AI outputs, which limited peer-to-peer explaining and design rationalizations. Gradually, this threatened to undermine the collaborative conversation that has long been a vehicle for transferring tacit architectural knowledge. Teams that compensated for this with formal protocols (e.g., requiring coders to restate the machine-authored rationale in pull requests or design documents) could most easily sustain stronger knowledge-sharing practices.

From a sociotechnical standpoint, developers experienced greater ease of use with AI support. Satisfaction results showed that participants felt that the AI partner alleviated the tedium of repetitive chores, enabling more cognitive room for solving higher-level problems. However, confidence in the accuracy of AI outputs was still inconclusive, and senior engineers expressed concern that we'd have to verify it constantly. This juxtaposition of desire for productivity versus a fear of inaccuracy reflects the bifurcated approach to AI adoption.

5. Discussion

The findings from this study emphasise the transformative nature, as well as the complexity of LLM-based copilots and agentic AI in influencing contemporary software engineering practices. Although we find evidence of large gains in productivity and some benefit to quality for the user firm, our empirical results also suggest that neither is automatic but rather depends on policies by the organization and its practices with respect to reviewing content, as well as the cultural fit of participation. The debate is therefore focused on the interpretation of these trade-offs and on ways towards sustainable integration.

The acceleration of regular tasks was the one benefit that all teams appreciated out of these. It does this by managing much of the generation and testing of boilerplate code, leaving developers to focus on higher-level problem-solving. These productivity gains, though, can be overestimated when considered in isolation. Factoring in downstream effects, such as defect leakage or increased review task load, the picture is dirtier. The findings indicate that productivity should not be gauged only by the speed of output, but rather as a sustained flow efficiency rolled up across the full life cycle. This difference is especially important for business, as speed improvements without quality guarantees that scale upwards can actually boost the total cost of ownership.

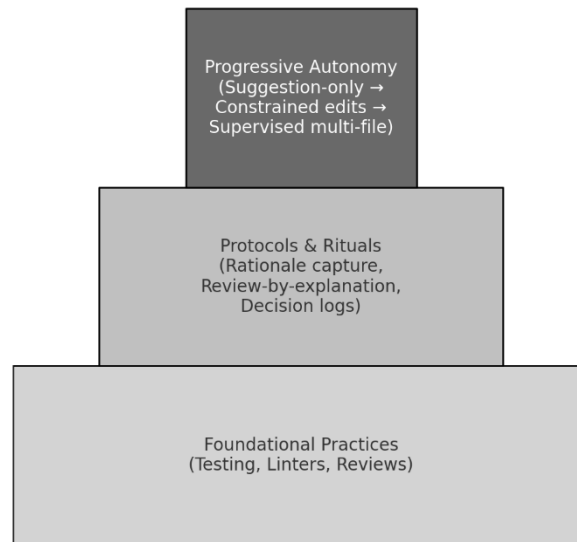


Figure 3: Governance Pyramid for Sustainable Human-AI Collaboration, Emphasizing Layered Safeguards from Engineering Practices to Rituals and Staged Autonomy

Quality results also reflect two realities. For one, the AI-authored code appears to adhere to stylistic conventions and pass basic validation checks more often than purely human-generated code. On the other hand, there are often mistakes due to nuances in meaning or context when we use AI suggestions without a critical eye. This highlights the importance of integrating AI within strong governance, where all recommendations receive the same validation pipeline as code written by a human. The results also highlight the need for vigilance on the part of reviewers. Teams that maintained strong review routines were able to turn AI-generated drafts into high-quality outputs, whereas those who cut reviews short to exploit the speed of AI saw more backsliding. So AI, far from relegating humans to the dustbin of history, instead increases the need for structured review processes as a last guard against substandard work.

The implications with respect to knowledge transfer are arguably the most far-reaching. The dynamics of the human-AI pair programming relationship are reframed, where the use of AI to expedite onboarding and offer scaffolding for junior developers also threatens to debilitate the dialogic reasoning upon which organisational learning depends. Over time, when developers blindly accept AI's outputs, the chances for peer-to-peer context generation lose ground and challenge/reflection are lost, which can wreak havoc on teams' long-term resilience. But, according to the research, this risk can also be countered with mindful behaviors. Practices like making developers justify AI-assisted changes in pull requests or keeping architectural decision records provide chances for reflective knowledge spreading. These found rituals of the technical and social life introduce AI4T as an engine for learning, not a replacement for it.

The findings also emphasize trust as a core issue in adoption. Developers, in the main, look kindly on AI partners, but they are less assured of their correctness – opinion on that has more to do with age and experience. This ambivalence is a specific example of a broader tension: there's enthusiasm about the productivity gains that can be realized and anxiety people have about whether they can trust the solutions produced. Hustling on trust will include technical advancements to ensure better model accuracy as well as organizational transparency on how AI tools are validated, monitored, and governed. Acceptance

rates, Change Ratios, and post-merge Defect Origins are metrics that can be used as feedback loops in order to tune the trust and foster responsible usage.

Last, this work is consistent with industry-wide skepticism about the adoption of agentic AI. Although the agents showed promising results for automation of the multi-step activities, they were susceptible to scope creep and context ambiguity. Without well-delimited boundaries and developmental autonomy models, agentic tools may generate rework rather than value. The best-performing deployments were the ones in which agents were limited to narrow, validation-gated tasks. It means enterprises need to do it in steps, starting with AI as a recommendation engine and moving on to limited-forced tool use before finally allowing supervised multi-file changes. This staged integration will make sure that the AI tools mature and grow alongside the socio-technical ecosystem.

6. Conclusion

Large language models and agentic AI systems integrated into software engineering teams are a paradigm shift in the way that we are building software. This work has shown how AI copilots and agentic tools yield observable productivity, contingent gains in code quality, and reconfigured knowledge transfer dynamics on teams. But the results also demonstrate that these advantages are not consistent, or necessarily automatic. A related set of concerns crystallizes most powerfully in disciplined contexts where the technical, organizational, and human governance align with opportunities as well as the limits posed by AI systems.

On the productivity side, the data is pretty clear that AI is great at speeding up tedious and well-bounded problems like code scaffolding, unit testing, or generated documentation. Cycle times were consistently faster, and copilots were used because they felt less cognitively fatiguing. These productivity improvements were further extended by agentic AI tools, which automated repository-level activities while combining multiple stages of workflows. But the findings also warn that productivity gains dwindle in projects that involve architectural reasoning or intricate trade-offs. In these scenarios, AI systems typically provide plausible answers that need heavy revision, reminding us that human expertise is still necessary for long-range decision making.

Code quality results echoed this complex picture. AI outputs scaled back unnecessary returns and got closer to sticking to the letter of stylistic and syntactic law, indicating a model that could internalize some sense of best practices. Pre-merge quality signals like test success rates and linting compliance grew in AI-assisted teams. However, continued subtle semantic defects, even when AI-generated submissions were accepted without back translation, underscored the importance of thorough review and validation. The research, therefore, serves as an important reminder that quality isn't assured by AI assistance, but rather a measure of how effectively organizations police their validation pipelines. Quality improved where disciplined review practices were maintained, and post-deployment regressions jumped where the review rituals were shortened.

The strategic ramification of your prioritizing how knowledge transfer can be conveyed is probably the most important. Human-AI teamwork changes the nature of team interaction itself. The process of getting new contributors onboarded became faster as copilots provided scaffolding at the touch of a button and explanations, which, with copy-paste bots, allowed total novices to make meaningful contributions sooner. However, diminished conversation depths in human-AI pair programming also caused alarm over the flattening of expertise transfer between peers and tacit information spreading. This finding indicates that, if left unaddressed, AI tools ultimately risk turning software teams into a more efficient but less resilient mechanical assembly line. Therefore, protocols like forcing reasons capture in pull requests, keeping architectural decision records, and incentives for developers to explain AI-derived reasoning are so important to keep the organizational knowledge alive.

Apart from the technical and team results, this research provides a broader social-technical rationale. Trust was identified as a facilitator and inhibitor of adoption. Developers appreciated the efficacy of AI assistance, but were wary of correctness, in particular for high-stakes or regulated domains. This dualism underscores the role of transparent governance, feedback loops, and phased adoption. Agency tools in particular need to be progressively rolled out with increasing autonomy: read-only/suggestion bus mode, constrained edits (with/ validation gates), and then be allowed for supervised multi-file activities. This maturity path also provides organizations with the reassurance that overreach risks are minimized, and they develop trust in the dependability of AI systems.

The larger point here is that human-AI collaboration is not just about enhancing personal productivity but eventually reshaping our socio-technical systems. Coking those organizations so narrowly on speed metrics can lead to some short-term "burst" productivity, but it's the definition of degrading your transitory competitive advantage and poisoning knowledge flows. On the other hand, those who invest in governance systems, cultural rituals, and transparent validation mechanisms will be able to make use of AI as an enabler for sustainable innovation. As such, the true value of AI to software teams isn't about replacing human ingenuity; it's about generating new forms of interaction between humans and machines that buttress one another.

This work should be expanded in the future to study long-term effects on architectural health, developer career paths, and cross-team coordination in multi-agent settings. Particularly, we should consider the psychological aspects of trust, dependence, and judgment when dealing with systems that are becoming more autonomous. As AI systems move from being copilots to collaborators, the lines separating suggestion, execution, and decision-making will further dissipate. Safeguarding these designs for human judgment and organizational goals will be the defining challenge of the next decade of software engineering.

References

1. S. Peng, et al., “The Impact of AI on Developer Productivity: Evidence from a Controlled Experiment with GitHub Copilot,” *arXiv preprint arXiv:2302.06527*, Microsoft Research, Feb. 2023.
2. GitHub, “Quantifying GitHub Copilot’s Impact in the Enterprise: Developer Productivity and Satisfaction,” GitHub White Paper, May 2024.
3. GitHub, “Research: Quantifying GitHub Copilot’s Impact on Developer Productivity and Happiness,” GitHub, Sept. 2022.
4. SWE-bench Team, “SWE-bench: A Benchmark for Repository-Level Software Engineering Tasks,” GitHub/Leaderboard Resources, 2023–2025.
5. W. Liang, Y. Zhang, and D. Lo, “An Exploratory Evaluation of Large Language Models Using Empirical Software Engineering Tasks,” in *Proc. Internetware Conf.*, ACM, July 2024.
6. C. E. Jimenez, et al., “Can Language Models Resolve Real-World GitHub Issues? Evaluating on SWE-bench,” *arXiv preprint arXiv:2310.06770*, Oct. 2023.
7. Y. Almeida, T. Pimentel, and A. Serebrenik, “AICodeReview: Advancing Code Quality with AI-Enhanced Code Review,” *Journal of Systems and Software*, vol. 209, 112922, Feb. 2024.
8. Stack Overflow, “2025 Developer Survey: AI Adoption and Trust in Software Engineering,” Stack Overflow Insights, July 2025.
9. A. Welter, M. Gerlach, and T. Fritz, “From Developer Pairs to AI Copilots: A Comparative Study on Knowledge Transfer in Pair Programming,” *arXiv preprint arXiv:2506.01234*, June 2025.
10. G. Fan, L. He, and X. Li, “Impact of AI-Assisted Pair Programming on Student Motivation, Anxiety, and Performance,” *International Journal of STEM Education*, vol. 12, no. 1, pp. 1–18, Jan. 2025.
11. S. Rasnayaka, H. Wang, and P. Liang, “An Empirical Study on Usage and Perceptions of Large Language Models in Software Engineering Projects,” *Empirical Software Engineering*, vol. 29, no. 3, pp. 1–28, Mar. 2024.
12. McKinsey & Company, “Unleashing Developer Productivity with Generative AI,” McKinsey Report, June 2023.
13. Reuters, “Over 40% of Agentic AI Projects Will Be Scrapped by 2027, Gartner Says,” Reuters Technology Report, June 25, 2025.
14. Microsoft Research, *The New Future of Work Report 2023: The Role of Generative AI in Knowledge Work*, Microsoft Research, Dec. 2023.
15. Microsoft Research, *Generative AI in Real-World Workplaces: Empirical Insights from Enterprise Deployments*, Microsoft Research, July 2024.