# The Impact of LLM-Assistants on Software Developer Productivity: A Systematic Literature Review

AMR MOHAMED, Queen's University, Canada

MARAM ASSI, Université du Québec à Montréal, Canada

MARIAM GUIZANI, Queen's University, Canada

Large language model assistants (LLM-assistants) present new opportunities to transform software development. Developers are increasingly adopting these tools across tasks, including coding, testing, debugging, documentation, and design. Yet, despite growing interest, there is no synthesis of how LLM-assistants affect software developer productivity. In this paper, we present a systematic literature review of 37 peer-reviewed studies published between January 2014 and December 2024 that examine this impact. Our analysis reveals that LLM-assistants offer both considerable benefits and critical risks. Commonly reported gains include minimized code search, accelerated development, and the automation of trivial and repetitive tasks. However, studies also highlight concerns around cognitive offloading, reduced team collaboration, and inconsistent effects on code quality. While the majority of studies (92%) adopt a multi-dimensional perspective by examining at least two SPACE dimensions, reflecting increased awareness of the complexity of developer productivity, only 14% extend beyond three dimensions, indicating substantial room for more integrated evaluations. Satisfaction, Performance, and Efficiency are the most frequently investigated dimensions, whereas Communication and Activity remain underexplored. Most studies are exploratory (64%) and methodologically diverse, but lack longitudinal and team-based evaluations. This review surfaces key research gaps and provides recommendations for future research and practice. All artifacts associated with this study are publicly available at https://zenodo.org/records/15788502.

## 1 INTRODUCTION

Large Language Models (LLMs) are increasingly being integrated into the software engineering (SE) domain [1]. In particular, the emergence of LLM-assistants, the term we use to refer to generative AI tools powered by LLMs that support software development tasks, has driven rapid adoption in both research and practice. Examples include OpenAI's GPT-series (e.g., GPT-4 [2]) and GitHub Copilot [3], which are now commonly used to assist with tasks such as code generation and completion[4, 5, 6], code translation [7, 8], debugging and maintenance [9, 10, 11], documentation [12, 13], and system design [14, 15]. These tools support a new development paradigm often referred to as AI pair programming, in which developers interactively engage with LLM-assistants throughout the software development

Authors' addresses: Amr Mohamed, amr.m@queensu.ca, Queen's University, Kingston, Ontario, Canada; Maram Assi, Université du Québec à Montréal, Québec, Montréal, Canada, assi.maram@uqam.ca; Mariam Guizani, mariam.guizani@queensu.ca, Queen's University, Kingston, Ontario, Canada.

process [16]. Since the public release of ChatGPT[1] in late 2022, an expanding ecosystem of LLM-powered coding assistants—such as Cursor[2], Windsurf[3], and Bolt[4] has emerged. This widespread integration underscores the growing reliance on LLMs in SE and raises critical questions about their impact on software developer productivity.

Software developer productivity is a multifaceted construct that encompasses not only the efficiency and quality of software production but also the satisfaction, collaboration, and cognitive load experienced by a developer. While early approaches for measuring productivity relied on quantifiable outputs such as lines of code (LOC) or development velocity [17, 18], recent research highlights the importance of human-centered factors such as communication, satisfaction, and well-being [19]. As LLM-assistants become increasingly integrated into development workflows, it is crucial to understand their impact on software developer productivity.

To address this gap, we conduct a systematic literature review (SLR) of 37 peer-reviewed studies published between 2014 and 2024 that examine the impact of LLM-assistants on software developer productivity. Our review analyzes methodological strategies, evaluation practices, and the productivity dimensions these studies engage with. We synthesize reported benefits and risks, and apply established conceptual frameworks to map our findings and contextualize their broader implications. Based on this synthesis, we identify key research gaps and provide actionable recommendations for both researchers and practitioners. This paper makes the following contributions:

- We present the first systematic literature review focused on the impact of LLM-assistants on software developer productivity, synthesizing evidence from 37 peer-reviewed primary studies published between 2014 and 2024.
- We provide a structured characterization of the methodological strategies and evaluation practices used to assess developer productivity, and synthesize the reported effects of LLM-assistants, surfacing key benefits (e.g., reduced task initiation overhead, support for code-adjacent tasks) and risks (e.g., over-reliance, flow disruption).
- We analyze our findings through the lens of the *SPACE* framework and employ McLuhan's Tetrad framework in our discussion to reflect on broader socio-technical implications.
- We offer actionable recommendations for practitioners and researchers, and release a publicly available replication package [20] containing all study data, selection decisions, and exclusion rationales to support transparency and reproducibility.

The remainder of this paper is structured as follows. Section 2 provides background on software developer productivity. Section 3 outlines the methodology used to conduct our SLR, including search strategies, selection criteria, and data extraction procedures. Sections 4 through 7 focus on addressing each of the research questions separately. Section 8 discusses implications, recommendations for practitioners, and directions for future research. Section 9 examines the threats to validity. Finally, section 10 concludes the paper.

## 2 BACKGROUND

Software developer productivity has received sustained research attention since the early decades of software engineering. Notably, Frederick J Brooks noted in his 1975 book *The Mythical Man Month* [21] that there is *"no silver bullet"* for improving productivity, underscoring that software development is a complex, human-centric activity not easily optimized by simply adding more resources. Since then, a wide range of studies have attempted to define and measure developer productivity [19, 22, 23, 24, 25].

---

[1] https://chatgpt.com/
[2] https://www.cursor.com/
[3] https://www.windsurf.com/
[4] https://bolt.new/

However, despite decades of investigations, developer productivity remains a complex concept, and a clear consensus on how to define or measure it accurately has yet to be established. The lack of a universally accepted definition or precise measurement consensus is largely due to the wide range of variables that influence productivity, including individual factors (e.g., developer fluency, motivation, experience), organizational practices (e.g., team collaboration, feedback), and contextual elements (e.g., task variety, remote work capabilities) [17]. Traditionally, software developer productivity has been defined as the ratio of output to input [25, 26, 27, 28, 29]. Outputs have been quantified using metrics such as lines of code added [22, 30], tasks completed [31], or functions implemented [32], while inputs have been measured in terms of time spent (e.g., hours, days, or months) [31, 33]. Although such metrics offer a quantifiable view of developer output, they often fail to capture the broader human, social, and organizational dimensions that characterize real-world software development.

Recent research has shifted toward multidimensional frameworks to more comprehensively assess software developer productivity. Noda et al. [34] propose a model based on Developer Experience (*DevEx*)[35], comprising three core dimensions: feedback loops, cognitive load, and flow state. Their work synthesizes insights from software engineering and human-computer interaction to operationalize how development environments and processes shape developers' day-to-day experiences. The *DevEx* model offers a practical framework for assessing and improving productivity from the developer's point of view. Forsgren et al. [19] introduce the *SPACE* framework, which characterizes software developer productivity across five dimensions: Satisfaction and well-being, Performance, Activity, Communication and collaboration, Efficiency and flow. Satisfaction and well-being refer to how fulfilled and healthy developers feel in relation to their work, tools, team, and organizational culture. Performance reflects the quality and effectiveness of the outcomes of a system or process, such as software reliability, absence of bugs, or customer satisfaction. Activity captures the count of observable work events or software artifacts such as commits, pull requests, code reviews, builds, or deployments. Communication and collaboration addresses how individuals and teams communicate, coordinate, and integrate their work through metrics like discoverability of documentation and expertise. Efficiency and flow focus on the uninterrupted progress of work at both individual and system levels. Consequently, the *SPACE* framework has been increasingly used in empirical studies to offer a more nuanced lens on productivity, especially in collaborative settings such as AI-assisted development [36, 37, 38, 39, 40, 41, 42].

The suggested multidimensional perspectives highlight that a single metric cannot meaningfully capture productivity and instead encourage the use of composite measures that reflect the complex and varied nature of software development work, including human-centric metrics.

## 3 SYSTEMATIC LITERATURE REVIEW METHODOLOGY

We aim to establish the current state of evidence on the effects of LLM-assistants on software developer productivity. To this end, we conduct a comprehensive analysis across three key dimensions: (1) the methodological strategies, procedures, and instruments employed in primary studies (2) the reported benefits and risks associated with the use of LLM-based assistance, and (3) the specific dimensions of developer productivity that have been investigated. Our overarching objective is to synthesize the fragmented body of existing knowledge, highlight methodological strategies and their instrumentation, and identify critical gaps to guide future research in this rapidly evolving area.

We ground our methodology in the seminal guidelines by Kitchenham and Charters [43], which are derived from evidence-based practices in medical research and have been adapted for use in SE.

In particular, this review is guided by the following research questions:

**RQ0: What are the characteristics of peer-reviewed studies that investigate the impact of LLM-assistants on software developer productivity?**

In RQ0, we contextualize the emerging research landscape surrounding LLM-assisted software development. Specifically, we examine this landscape from three angles: (1) an overview of the temporal distribution of publications, (2) the publication venues and their disciplinary focus, and (3) the patterns of authorship across the research community.

**RQ1: What are the methodological strategies, procedures, and instruments used by peer-reviewed studies that investigate the impact of LLM-assistants on software developer productivity?**

In RQ1, we examine how existing research is conducted to investigate the relationship between LLM-assistants and software developer productivity. Specifically, we investigate the empirical strategies adopted to study the impact of LLM-assistants, (2) the procedures and study designs used to carry out these investigations, and (3) the instruments and metrics employed to evaluate software developer productivity.

**RQ2: What is the impact of LLM-assistants on software developer productivity?**

In RQ2, we explore the effects of LLM-assistants on software development practices. First, we summarize the overall findings from the identified primary studies. We supplement these findings by analyzing the reported benefits and risks of using LLM-assistants across diverse study settings. This question aims to provide a structured understanding of how these tools affect developers in practice and what trade-offs they introduce.

**RQ3: Which dimensions of developer productivity are investigated and how do these dimensions map onto the SPACE framework?**

In RQ3, we investigate how productivity is defined and assessed in the context of LLM-assisted software development. To provide additional insights, we map the main focus of each study to the dimensions of the *SPACE* framework, i.e., Satisfaction and well-being, Performance, Activity, Communication and collaboration, and Efficiency and flow. This analysis offers a clear understanding of how the concept of productivity is operationalized across the existing body of work and highlights underexplored dimensions to guide future research.

### 3.1 Pre-review mapping

We adhere to the guidelines proposed by Kitchenham and Charters [43] for piloting the research protocol through pre-review mapping. To plan the construction of our review, we conducted a pre-review planning study, as outlined in Figure 1. Piloting the research protocol spanned a three-month period.

*3.1.1 Control papers identification.* Once our research questions defined, we set our list of inclusion and exclusion criteria (see section 3.1.1). To identify control papers, we manually search Google Scholar using various keywords including "Productivity," "Large Language Models," "Software Developer," and "ChatGPT". After screening the title and abstract for relevance, we obtain an initial set of 20 candidate papers. We apply one iteration of forward and backward snowballing to expand the candidate set. Specifically, we examine the reference list in each candidate article (i.e., backward snowballing) and identify studies that cite them (i.e., forward snowballing). These two steps results in 15 additional articles.

Table 1. Database search strings and results. Total n = 8,540.

| Database | Search String | Results (since 2014) |
|---|---|---|
| ACM | (Language Model* OR "LM" OR "LMs" OR "LLM" OR "LLMs" OR "Artificial Intelligence" OR "AI") AND ((title:(Software Engineer* OR Software Develop* OR Developer* OR Coder* OR Programmer*)) OR (abstract: (Software Engineer* OR Software Develop* OR Developer* OR Coder* OR Programmer*))) AND (Productivity) | 4,044 |
| IEEE Xplore | (((Language Model OR "LM" OR "LMs" OR "LLM" OR "LLMs" OR "Artificial Intelligence" OR "AI") AND (("Document Title":Software Engineer OR "Document Title":Software Develop* OR "Document Title":Developer OR "Document Title":Coder OR "Document Title":Programmer) NEAR/5 (Productivity)) OR (("Abstract":Software Engineer OR "Abstract":Software Develop* OR "Abstract":Developer OR "Abstract":Coder OR "Abstract":Programmer) NEAR/5 (Productivity)))) | 491 |
| ScienceDirect | ((Language Model OR "LM" OR "LMs" OR "LLM" OR "LLMs" OR "Artificial Intelligence" OR "AI") AND (Productivity)) <br> Advanced Search: Title, abstract, keywords: (Software Engineer OR Software Development OR Developer OR Coder OR Programmer) | 3,734 |
| Web of Science | ALL=(Language Model OR "LM" OR "LMs" OR "LLM" OR "LLMs" OR "Artificial Intelligence" OR "AI") AND ((TI=(Software Engineer NEAR Productivity OR Software Develop* NEAR Productivity OR Developer NEAR Productivity OR Coder NEAR Productivity OR Programmer NEAR Productivity)) OR (AB=(Software Engineer NEAR Productivity OR Software Develop* NEAR Productivity OR Developer NEAR Productivity OR Coder NEAR Productivity OR Programmer NEAR Productivity))) | 271 |

All authors independently conduct a full-text screening and assess all retrieved articles against the predefined inclusion criteria. Three meetings were held to discuss and resolve any disagreements. This process results in a final set of 17 articles, which we use to guide the construction and testing of our search string. The set of control papers can be found in the supplemental material [20].

**Inclusion and Exclusion Criteria.** We select papers for our study based on the following inclusion (IC) and exclusion (EC) criteria. The supplemental material [20] includes the exclusion decision for each study based on specific IC & EC criteria.

**Inclusion (IC)**:

- (+) IC1: The paper investigates the effect of AI or LLMs on software developer productivity.
- (+) IC2: The paper is in English.
- (+) IC3: The paper has an accessible full text and was published in 2014 or later.

**Exclusion (EC)**:

- (−) EC1 (Out of Scope): The paper does not focus on SE or does not explore the impact of AI or LLMs on software developer productivity.
- (−) EC2 (Out of Focus): The paper mentions the impact of AI or LLM on software developer productivity without it being one of the topics of the study.
- (−) EC3 (Publication Type): The paper belongs to any of the following categories: secondary studies; work-in-progress, extended abstracts, posters, tool demos, editorials, or grey literature; studies published in books,
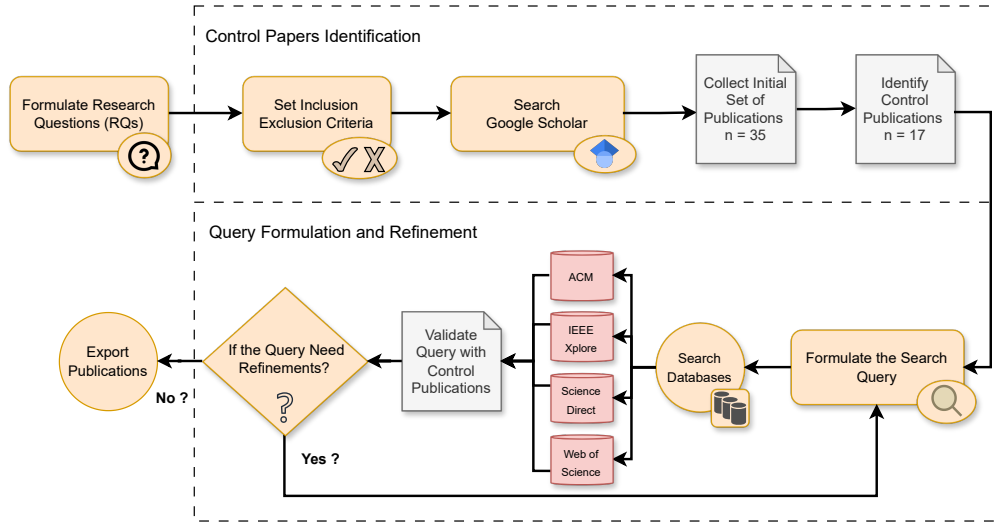
Fig. 1. Overview of our pre-review mapping, including the identification of control papers, and query formulation and refinement.

theses, workshop, monographs, keynotes, panels, doctoral symposium, or any other venues without a formal peer-review process.

- (−) EC4 (Length): The paper is a short publication with fewer than four pages.

*3.1.2 Query formulation and refinement.* Following the guidelines proposed by Kitchenham and Charters [43], we select four major digital libraries that are widely used in SE research: ACM Digital Library[5], IEEE Xplore[6], ScienceDirect[7], and Web of Science[8].

Constructing effective search queries is a challenging step in developing an SLR protocol, particularly due to the absence of standardized guidelines for selecting search terms [43]. This process relies on iterative refinement [44], especially in contexts where terminologies such as those used to describe AI and LLM in SE vary significantly between studies. Given the growing volume of research in this area, we observe that overly broad queries return a large number of false positives, thereby reducing the precision of the search query. Conversely, narrow queries risk omitting relevant studies [45].

The first and last authors held five meetings to iteratively develop and refine the search string. In each iteration, we test candidate queries by checking whether they successfully retrieve all 17 control studies identified during the pilot phase [45] (see Figure 1). Based on these tests, we adjust both the keywords and the structure of the queries to improve their relevance and coverage. For example, we refine the search strings by replacing terms like "software engineering" and "software engineers" with a wildcard term "software engineer*", following the syntax and capabilities of each database. We also avoid using the keywords "assessment", "performance", and "efficiency" after observing a high volume of irrelevant studies that focus on technical-level benchmarks rather than productivity. After five query
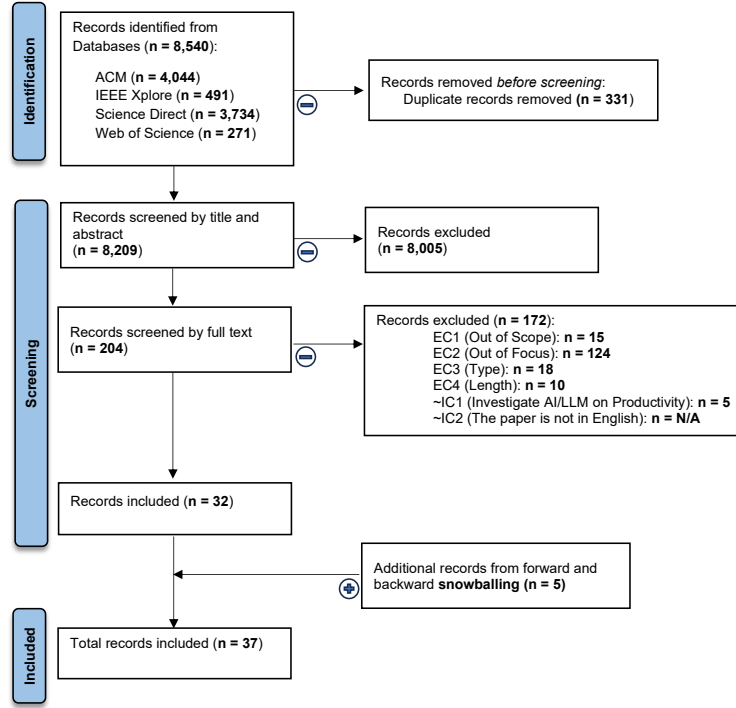
---

Fig. 2. Overview of the selection process for primary studies included in the SLR using PRISMA flow chart [49].

iterations, all authors held a consensus meeting and agreed on the final search string (see Table 1). The final selected search query successfully retrieves all 17 control papers.

Table 1 presents the finalized search query used in our review. Each search query consists of three segments separated by the AND string. The first segment related to AI or LLMs limits the filter to the AI or LLM technology. The second segment refers to the actor (i.e., software developers), and the third segment captures the concept in question (i.e., productivity).

In alignment with recent SLRs [46, 47, 48], we restrict our searches to the title, abstract, and keywords, as this improves precision. We also leverage proximity operators in the query, namely "NEAR/5", to improve the contextual relevance of matched terms [26]. This operator retrieves documents where the specified terms appear within five words of each other. This refinement is only applied in IEEE Xplore and Web of Science, as the ACM Digital Library and ScienceDirect do not support it. All searches were conducted on December 31, 2024.

## 3.2 Primary Study Selection Process

We execute the search protocol as illustrated in Figure 2. First, we apply the finalized search strings on the selected digital libraries, restricting the results to publications from 2014 onward. This yields an initial set of (n = 8,540) records, including 4,044 papers from ACM, 491 papers from IEEE Xplore, 3,734 papers from ScienceDirect, and 271 papers from
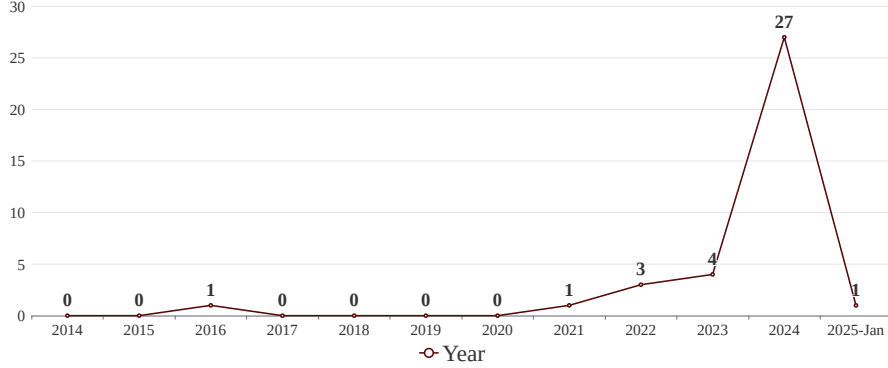
Fig. 3. Publication frequency per year.

Web of Science. After removing duplicates, we obtain a set of 8,209 unique records (n = 331 excluded). The first author performs an initial screening of all records based on their titles and abstracts, a process that took a total of 41 days. We use the tool *Rayyan*[9] to tag any excluded paper with the corresponding exclusion criteria (the list of all excluded papers and their corresponding exclusion criteria is provided in the supplemental material [20]). The second and last authors independently validate the excluded papers. Three meetings were held to discuss any disagreements until reaching a consensus. We adopt a conservative screening approach whereby records with insufficient information in the title and abstract to support a clear inclusion or exclusion decision were included in full-text review. The title and abstract screening process resulted in 204 papers (see Figure 2), from which we exclude a total of 8,005 records.

The full-text screening phase was conducted over a period of 7 weeks. This phase involves the careful reading and evaluation of the remaining 204 papers, applying the inclusion and exclusion criteria to assess their eligibility. For a detailed rationale on the inclusion and exclusion criteria during full-text screening, please refer to the supplemental material [20]. Throughout this process, the first author led the full-text screening of papers. Whenever the relevance of a study was unclear, it was flagged and reviewed in consultation with the second and last authors to ensure alignment with the protocol. This results in the exclusion of 172 studies during the full-text screening process.

To further expand our study set, we conduct a snowballing procedure by examining both the references and citations of the 32 selected studies. This phase took approximately two weeks, resulting in the identification and inclusion of five additional articles, expanding the final set to 37 primary studies. The thematic analysis and synthesis of qualitative findings for RQs across these studies were conducted over a period of three months.

## 4 RQ0: WHAT ARE THE CHARACTERISTICS OF PEER-REVIEWED STUDIES THAT INVESTIGATE THE IMPACT OF LLM-ASSISTANTS ON SOFTWARE DEVELOPER PRODUCTIVITY?

### 4.1 Publication years

Figure 3 shows the frequency of publications by year. Although we included studies published within the ten years preceding our review (2014–2024), only four were published between 2014 and 2022 [PS1, PS2, PS3, PS4]. Research interest began to rise in 2022, which coincides with the release of ChatGPT[10]. This culminated in a sharp peak in 2024,

---

[9]https://rayyan.ai/
[10]ChatGPT: https://openai.com/chatgpt

which accounts for 73% of all included studies, likely due to increased accessibility and interest following ChatGPT's release.

## 4.2 Author distributions

We investigate the distribution of papers per author across the 143 authors of primary studies. Most authors (i.e., 135) have a single publication, while 7 authors have two papers, and the most prolific author, Igor Steinmacher, has three publications. This distribution is probably due to the fact that the investigation of the impact of LLM-assistants on software developer is a relatively new topic that is just starting to build momentum.

Table 2. Distribution of primary studies by publication venues.

| Research Focus | Venue | Primary Studies | % |
|---|---|---|---|
| Software Engineering and Computer Science | ACM Transactions on Software Engineering and Methodology (TOSEM) | [PS2, PS5, PS6, PS7, PS8] | 54% |
| | International Conference on Software Engineering (ICSE) | [PS9, PS10, PS11] | |
| | Proceedings of the ACM on Software Engineering (PACMSE) | [PS12, PS13, PS14, PS15] | |
| | ACM SIGPLAN International Symposium on Machine Programming (PLDI) | [PS16] | |
| | Proceedings of the ACM SIGSOFT International Symposium on Foundations of Software Engineering (FSE) | [PS1] | |
| | Software Engineering in Practice (ICSE-SEIP) | [PS17] | |
| | Automated Software Engineering (ASE) | [PS18] | |
| | Software Quality, Reliability, and Security Companion (QRS-C) | [PS19] | |
| | Science of Computer Programming | [PS20] | |
| | Conference on the Art, Science, and Engineering of Programming | [PS21] | |
| | Evaluation and Assessment in Software Engineering (EASE) | [PS22] | |
| Human-Computer Interaction (HCI) | ACM Conference on Human Factors in Computing Systems (CHI) | [PS4, PS23] | 16% |
| | Proceedings of the ACM on Human-Computer Interaction (CSCW) | [PS24] | |
| | ACM Transactions on Computer-Human Interaction (TOCHI) | [PS25] | |
| | International Conference on Intelligent User Interfaces (IUI) | [PS3, PS26] | |
| AI for Software Engineering / AI Engineering | AI Engineering - Software Engineering for AI (IEEE/ACM) | [PS27] | 8% |
| | AI-Powered Software (AIware) | [PS28] | |
| | International Conference on Generative Artificial Intelligence and Information Security (GAIIS) | [PS29] | |
| Information Systems and Decision Science | Journal of Decision Systems (JDS) | [PS30] | 8% |
| | Proceedings of the Americas Conference on Information Systems (AMCIS) | [PS31] | |
| | Innovations in Software Engineering Conference (ISEC) | [PS32] | |
| Human-Aspects and Socio-Economic Impact | Futures | [PS33] | 8% |
| | Structural Change and Economic Dynamics | [PS34] | |
| | Cooperative and Human Aspects of Software Engineering (CHASE) | [PS35] | |
| Software Engineering Education | Innovation and Technology in Computer Science Education (ITiCSE) | [PS36] | 5% |
| | ICSE Software Engineering Education and Training (ICSE-SEET) | [PS37] | |

### 4.3   Publication venues

We extract the publication venue for each primary study. Table 2 shows the venues categorized by research focus. The majority (54%) of primary studies fall under"Software Engineering and Computer Science" published in venues including TOSEM, ICSE, and EASE. Human-Computer Interaction (HCI) is the second most prominent research focus with 16% (6 out of 37) primary studies published in venues such as CHI and TOCHI. The remainder of the primary studies are similarly distributed among specialized venues focusing on AI for Software Engineering / AI Engineering, Software Engineering Education, Information Systems and Decision Science, and Human-Aspects and Socio-Economic Impact. The variety of publication venues highlights the breadth and depth of the topic under study. The integration of LLM-assistants into the software development workflow introduces important considerations related to usability, automation, interaction design, and developer behavior.

> **RQ0- Summary**
>
> The majority of studies (33 out of 37) have been published after the release of ChatGPT in November 2022. Only four earlier studies (2014–2022) use pre-LLM paradigms. The majority of authors (135 out of 143) have only a single publication included, and 7 authors have two papers. Most publications appear in Software Engineering and Computer Science venues (20 studies), followed by HCI (6 studies).

## 5   RQ1: WHAT ARE THE METHODOLOGICAL STRATEGIES, PROCEDURES, AND INSTRUMENTS USED BY PEER-REVIEWED STUDIES THAT INVESTIGATE THE IMPACT OF LLM-ASSISTANTS ON SOFTWARE DEVELOPER PRODUCTIVITY?

### 5.1   Distribution of the research strategies

Table 3. Distribution of research strategies across the primary studies.

| Strategy | Primary Study | Percent |
|---|---|---|
| Field Study | [PS12, PS14, PS17, PS18, PS28, PS29, PS31, PS35] | 22% |
| Field Experiment | [PS37] | 3% |
| Experimental Simulation | [PS20, PS23, PS32] | 8% |
| Laboratory Experiment | [PS2, PS3, PS4, PS7, PS10, PS11, PS13, PS15, PS19, PS22, PS24, PS25, PS26, PS27, PS36] | 41% |
| Sample Study | [PS6, PS9, PS16, PS30, PS34] | 14% |
| Judgment Study | [PS33] | 3% |
| Formal Theory | [PS1, PS8] | 5% |
| Others (e.g., opinion papers) | [PS5, PS21] | 5% |

We classify the 37 primary studies based on Stol and Fitzgerald [50] taxonomy for empirical software engineering strategies (see Table 3), which is a taxonomy built to distinguish studies' strategies according to their level of obtrusiveness (e.g., control) and generalizability (e.g., realism). This classification offers a structured understanding of how the research community has approached the investigation of AI tools. Laboratory experiments are the most common strategy, used by 41% of the primary studies (15 out of 37). Laboratory experiments rely on controlled environments to isolate the effects of LLM-assistants on specific development tasks. Field studies are the second most common strategy, representing 22% (8 out of 37) of the primary studies. Field studies prioritize ecological validity by observing developer behavior in real-world settings without the need for researcher intervention. Sample studies, typically large-scale
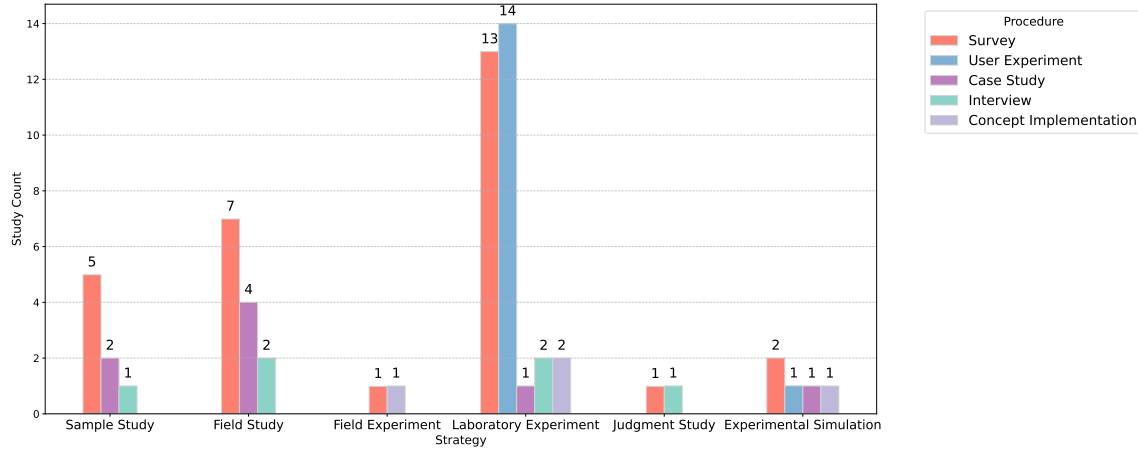
Fig. 4. The distribution of empirical procedures across methodological strategies.

surveys, account for 14% (5 out of 37), they aim to capture broad trends and perceptions across diverse developer populations.

Other strategies are less frequent. Field experiments (3%, 1 out 37) resemble field studies in that they occur in real-world settings. However, unlike field studies, researchers actively manipulate specific variables such as changes to a tool or process to evaluate their effects within the natural context. Experimental simulations (8%, 3 out of 37) combine elements from both laboratory and field studies by replicating real-world scenarios within controlled environments to study certain phenomena under realistic but simulated conditions. Judgment studies (3%, 1 out of 37) involve collecting expert opinions in a structured manner in a series of interviews and questionnaires.

The vast majority 89% (33 out of 37) of the primary studies employ an empirical research strategy, and only four studies are non-empirical. Given that this research area is relatively new, such studies may offer valuable context and early insights that can guide and support future empirical investigations.

## 5.2 Distribution of the research procedures

Table 4. Methodological procedures used in the empirical studies.

| Procedure | Primary Studies | Percent |
|---|---|---|
| Survey | [PS2, PS3, PS4, PS6, PS9, PS10, PS11, PS12, PS14, PS15, PS16, PS17, PS18, PS19, PS20, PS22, PS23, PS24, PS25, PS26, PS27, PS28, PS29, PS30, PS31, PS33, PS34, PS36, PS37] | 88% |
| User Experiment | [PS2, PS3, PS4, PS10, PS11, PS13, PS15, PS19, PS22, PS23, PS24, PS25, PS26, PS27, PS36] | 45% |
| Case Study | [PS3, PS16, PS17, PS28, PS29, PS30, PS31, PS32] | 24% |
| Interview | [PS4, PS25, PS29, PS30, PS33, PS35] | 18% |
| Concept Implementation (Proof of Concept) | [PS2, PS7, PS32, PS37] | 12% |

We adapt a fine-grained taxonomy identified from prior literature Glass, Vessey, and Ramesh [51] to assign one or more methodology procedures to each primary study, such as interviews, surveys, or user experiments (i.e., controlled experiments, quasi-experiments), based on the primary methods used (see Table 4).
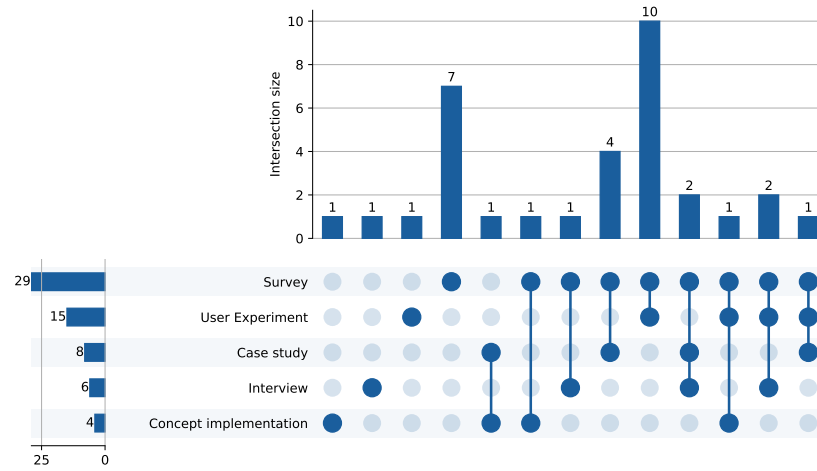
Fig. 5. Classification of methodological procedures classification and their overlap. Each bar on the left represents the number of studies that use a specific research method, while the bars along the top indicate how many studies employ a given combination of methods. For example, the most common pairing is a user experiment combined with a survey, used in 10 unique studies.

Among the empirical primary studies, we find 91% (30 out 33) of the studies leverage self-reported data such as surveys and interviews, and 45% of the studies (15 out of 33) are conducted in experimental settings (e.g., user studies, controlled experiments, or quasi-experiments). User experiments are almost exclusively associated with the laboratory experiment research strategy as shown in Figure 4. Additionally, 67% of these studies (22 out of 33) adopt a mixed-method approach as illustrated in Figure 5). For example, the most common combination of methods is a "user experiment" paired with a "survey". This approach is frequently used to triangulate self-reported perceptions (e.g., user experience or satisfaction) with measured performance metrics.

To provide additional insight into the nature of the current research direction, we assess the empirical studies based on their objectives, adopting a taxonomy from Hartson et al. [52] to classify the objective of each study into one of two approaches: formative and summative (see supplemental material [20] for detailed classification). A study with a formative objective primarily focuses on exploring, refining, or improving a process, tool, or methodology. In contrast, a study with a summative objective focuses on drawing conclusions about the effectiveness, outcomes, or impact of a completed process, tool, or methodology. We find that 64% (21 out of 33) of the studies have a formative goal, and 36% (12 out of 33) have a summative goal. This demonstrates that the current research direction is primarily exploratory, with a strong emphasis on early-stage development and evaluation rather than on drawing a final, conclusive outcome.

We also classify each empirical study based on its adopted methods of data analysis: quantitative, qualitative, or both (see supplemental material [20] for classification details). 73% (24 out of 33) of the empirical primary studies include a mix of quantitative and qualitative analysis, 12% (4 out of 33) of the studies rely only on qualitative analysis, and 15% (5 out of 33) of the studies rely only on quantitative analysis.

Table 5. Data sources and origin used by empirical studies.

| Data Source | Instrument Origin | Instrument and Primary Studies |
|---|---|---|
| **Self-Reported** | Designed by Authors | Surveys [PS2, PS3, PS6, PS9, PS14, PS15, PS17, PS19, PS20, PS22, PS24, PS26, PS28, PS29, PS33, PS37] |
| | | Interviews [PS4, PS25, PS29, PS30, PS33, PS35] |
| | | Users open-ended feedback [PS12, PS18, PS25, PS27] |
| | Validated Instruments and Frameworks | NASA-TLX (Mental Effort) [PS3, PS10, PS11, PS15, PS25, PS36] |
| | | Technology Acceptance Model (TAM) [PS6, PS11, PS37] |
| | | SPACE Framework-Based Surveys [PS16, PS23, PS24, PS31] |
| | | Self-Efficacy Questionnaires [PS4, PS10] |
| | | After-Action Review for AI (AAR/AI) [PS10] |
| | | Emotion Affect Questionnaire [PS36] |
| **Behavioral & Performance Metrics** | Designed by Authors | Task Completion and Correctness [PS2, PS10, PS11, PS13, PS26, PS36] |
| | | Suggestions Acceptance Rate [PS7, PS12, PS13, PS16, PS18, PS23, PS24] |
| | | Interaction Patterns (Logs/Edits/Tracking) [PS2, PS11, PS15, PS18, PS24, PS25, PS26] |
| | | Time to Completion [PS2, PS4, PS11, PS13, PS15, PS19, PS24, PS25, PS26, PS31, PS32] |
| | | Code Quality Metrics [PS2, PS3, PS10, PS19, PS24, PS29] |
| | | Productivity Gain [PS32] |
| | Validated Frameworks | Time Cost Quality (TCQ) Framework [PS34] |
| | | Resource-Based View (RBV) Framework [PS30, PS34] |

## 5.3 Evaluation instruments

Researchers employ various instruments, from self-reported surveys and interviews including validated questionnaires to behavioral & performance metrics as shown in Table 5. Self-reported methods remain predominant, often designed by study authors to capture user experience, perceived productivity, trust, or ease of use (e.g., post-task surveys, open-ended feedback). Only a subset of the empirical studies (15 out of 33) incorporate validated instruments including the *SPACE* framework [19], NASA-TLX for mental workload [53], TAM for technology acceptance [54], self-efficacy questionnaires [55, 56], and emotional affect questionnaire [57]. Behavioral & performance metrics focus on quantifiable outcomes, such as time to completion, acceptance rate of AI-generated suggestions, code quality metrics, and some analysis of interaction patterns (see Table 5). We find that behavioral & performance metrics are mostly associated with studies with a high level of control, such as laboratory experiments, field experiments, or experimental simulation (see Figure 6), for example, metrics such as time to completion or code quality metrics are mainly associated with laboratory experiments. In contrast, field studies and sample studies, while still employing diverse sets of instruments, primarily rely on self-reported methods, such as surveys, interviews, and users' open-ended feedback.

*5.3.1 Time to completion.* Measuring the time required to complete tasks is the most frequently used performance metric. 30% (11 out of 33) of the empirical primary studies employ this measure. The majority of studies that measured task completion time are laboratory experiments [PS2, PS4, PS11, PS13, PS15, PS19, PS24, PS25, PS26], which assess the time taken to complete specific programming tasks under controlled conditions. One field study conducted within a software company [PS31] measures time-related performance by comparing throughput and cycle time before and after the integration of Copilot. Additionally, an experimental simulation [PS32] estimates total effort using a "person/days" metric to compare task completion durations with and without LLM-assistants.
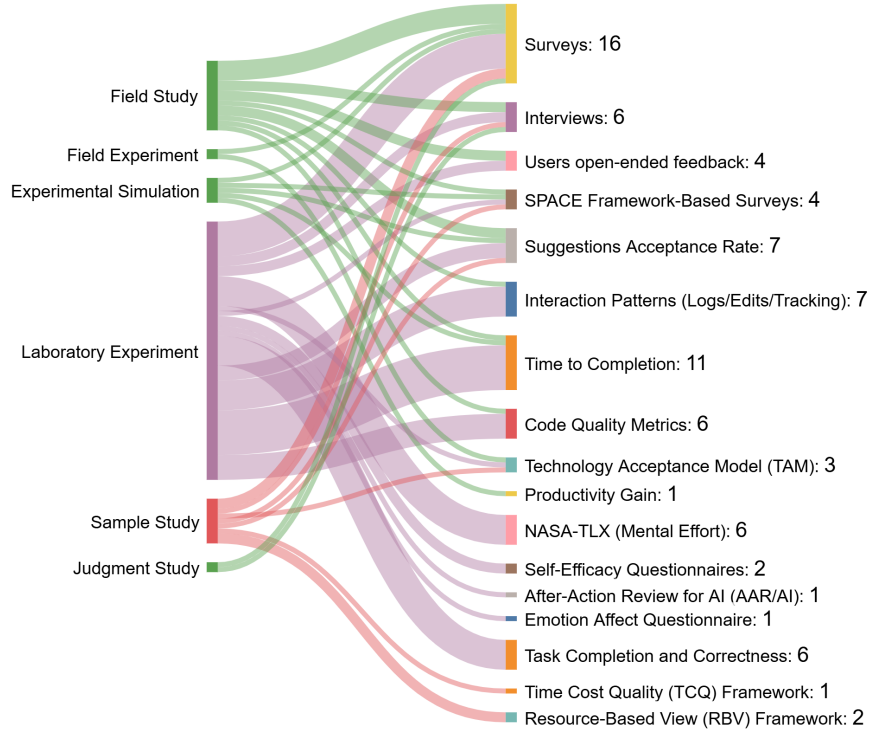
Fig. 6. Mapping between empirical strategies and instruments in primary studies.

*5.3.2 LLM suggestions acceptance rate.* Measuring the acceptance rate of LLM suggestions is one of the commonly used behavioral and performance metrics to measure productivity [PS7, PS12, PS13, PS16, PS18, PS23, PS24]. For instance, a study conducted by Meta [PS12] evaluates the adoption of an internal LLM-assistant coding system (Code Compose) within the company. The study quantifies the LLM-assistant's utility by measuring both the number of LLM-generated suggestions accepted by developers and the proportion of code authored by the LLM-assistants. These metrics are then compared against reported acceptance rates from competing LLM-assistants to evaluate relative effectiveness.

One reason the acceptance rate metric is widely adopted is a study conducted by GitHub [PS16], which statistically analyzes the relationship between several interaction metrics related to code completion and developers' self-reported productivity. The findings reveal a strong correlation between the frequency of accepted suggestions and perceived productivity. Despite these findings, the authors caution against using this metric in isolation to assess the effectiveness of LLM-assistants. They highlight that optimizing for acceptance rate may bias LLM-assistants toward well-represented languages or routine tasks, potentially disadvantaging less-represented workflows. Moreover, they warn that "blind" reliance on acceptance rate can lead to superficial improvements that inflate perceived usefulness without meaningfully enhancing developer outcomes.

*5.3.3 Mental effort and cognitive load.* Studies often use the terms mental effort or cognitive load interchangeably [58]. Reducing mental effort is considered a motivation for incorporating LLM-assistants in software development. Traditionally, studies aim to measure developer cognitive load through biometric modalities, such as electroencephalogram (EEG) or electrocardiogram (ECG) [58]. None of the identified studies leverages any biomedical measures or sensors for measuring cognitive load. Only one experimental study leverages eye-tracking [PS13] to measure the time participants spend reading code documentation.

Six studies (6 out of 33) [PS3, PS10, PS11, PS15, PS25, PS36] aim to measure cognitive workload by using NASA-TLX [53], which is a widely used questionnaire for assessing perceived mental workload. It captures six dimensions: mental demand, physical demand, temporal demand, performance, effort, and frustration. All six studies measure cognitive load in comparative experimental settings. We identify mixed findings regarding LLM-assistants' impact on mental cognitive load. In fact, a set of studies reports improvements [PS11, PS26, PS36], others neutral effects [PS3, PS15], and only one study reports a significantly worse experience in terms of frustration level [PS10]. For instance, [PS36] reports that Copilot reduces both perceived effort and mental demand for novice programmers. Similarly, [PS26] develops a custom questionnaire to assess cognitive load during programming exam tasks. Their findings show that students using Google Bard report lower mental effort compared to those relying on conventional search engines.

In contrast, [PS10] observes no significant difference in overall cognitive load between students using ChatGPT (GPT-4) and those using a traditional web browser, but does report a statistically significant increase in frustration for the ChatGPT group. Similarly, [PS3] finds that participants rate LLM-assisted tasks as equally demanding and effortful as tasks completed without LLM-assistants. Lastly, [PS15] finds no statistically significant differences across all NASA-TLX dimensions when comparing coding with and without ChatGPT (GPT-3.5).

The variability in reported effects of LLM-assistants on cognitive load highlights the complexity of evaluating mental effort in software development settings. These differences likely stem from diverse operationalizations of cognitive load, differences in participants' expertise, task design, and the capabilities of LLM-assistants across studies. This highlights the need for more standardized methodologies and multi-modal assessment strategies to draw robust conclusions about the cognitive impact of LLM-assistants.

*5.3.4 Econometric analysis.* Productivity is a concept primarily inherited from economics and project management. Two complementary studies investigate the impact of LLM-assistants using quantitative econometric analysis of productivity metrics [PS30, PS34] (see Table 5). The first study [PS34] leverages Time-Cost-Quality (TCQ) conceptual framework to conduct a comparative survey of over 1,000 large firms from 2021 to 2023, examining the effect of GenAI on labor productivity across different domains, including coding and content production.

Productivity is assessed in terms of both throughput (i.e., time efficiency) and quality (i.e., correctness of output). The study finds that coding exhibits the highest reported gains, with an average 24% improvement in throughput and 26% in quality. A complementary study by the same author [PS30] investigates the use of LLM-based pair programming through a survey of 70 large global companies. While the findings confirm that LLM-assistants can enhance development throughput, the study also identifies a critical trade-off: increased throughput is negatively correlated with code quality (r = −0.45). The study suggests that while LLM-assistants can enhance productivity, their effectiveness depends heavily on organizational readiness and the ability to balance speed with software quality.

---

**RQ1: Summary**

Empirical research dominates this space (89%), with laboratory experiments as the most common strategy (41%). Mixed-methods designs are prevalent (67%) among empirical methods, often combining user experiments with surveys. Time to completion is the most frequently used performance metric. Acceptance rate is a frequently used behavioral metric, though some studies caution against its overuse. Cognitive load findings are mixed: 5 studies use NASA-TLX, but results vary from reduced effort to increased frustration. Overall, the research is exploratory (64% formative) and methodologically diverse, highlighting early-stage evaluation and trade-offs between productivity and quality.

---

## 6  RQ2: WHAT IS THE IMPACT OF LLM-ASSISTANTS ON SOFTWARE DEVELOPER PRODUCTIVITY?

We conduct a thematic analysis of the findings reported in each primary study. Our analysis reveals several recurring themes, with multiple studies identifying common benefits and risks associated with the use of LLM-assistants. Figure 7 summarizes the frequency of discovered themes, where each theme is reported as a benefit or risk across the primary studies.
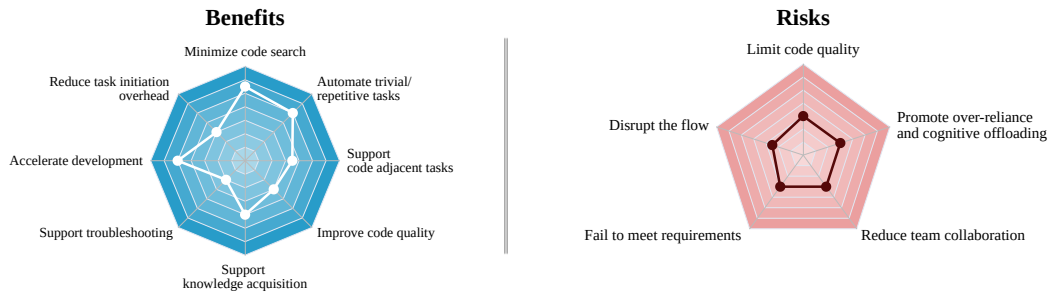


Fig. 7. Radar plots summarizing how frequently each theme appears as benefit (left) and risk (right) across primary studies on LLM-assisted development, scaled to a maximum of 14 studies, with a step size of two.

### 6.1  Benefits

*6.1.1  Minimize online code search.* Minimizing online code search is the most frequently discovered finding that highlights LLMs benefit (see Figure 7). Indeed, multiple studies highlight the potential of LLM-assistants to reduce the effort required for information retrieval [PS6, PS9, PS12, PS17, PS23, PS28, PS35]. For example, participants in [PS35] report that these tools are particularly advantageous for code-snippet retrieval compared to traditional Q&A platforms such as Stack Overflow.

Reducing the need for online search offers several benefits, including helping developers maintain a state of flow (e.g., "stay in the flow") [PS23], improving the speed of syntax recall [PS9], facilitating the discovery of unfamiliar APIs [PS9, PS12], and providing an alternative in situations where online search methods fail to deliver [PS6]. LLM-assistants are perceived as a more productive alternative to conventional code search, even when additional effort is required for validation [PS28]. Several controlled experiments show the benefits of using LLM-assistants compared to traditional online search [PS2, PS11, PS24]. In a user study involving 24 participants, [PS24] employs the *SPACE* framework to compare three conditions: traditional web search, code completion (i.e., Copilot), and interaction with a conversational

Table 6. Summary of LLM-assistants benefits on developer productivity

| Theme | Summary |
|---|---|
| Minimize online code search | Study participants in [PS6, PS9, PS12, PS17, PS23, PS28, PS35] noted that LLM-assistants reduce the effort of traditional online search. Controlled experiments also report benefits over online search [PS11, PS15] with some studies having mixed results [PS2, PS19]. |
| Accelerate software development | Studies highlight through self-reported methods that LLM-assistants accelerate software development [PS9, PS12, PS20, PS23, PS27, PS35], and quantitative measures demonstrate that LLM-assistants can reduce task completion time [PS11, PS15, PS25, PS32]. |
| Automate trivial/ repetitive tasks | LLM-assistants help minimize repetitive coding [PS6, PS9, PS14, PS33] by generating boilerplate code [PS3, PS9, PS12, PS20, PS31, PS35] and reducing keystrokes and typing effort [PS9, PS17]. A Delphi study [PS33] highlighted a future scenario where routine tasks are fully automated, freeing developers to focus on more complex work. |
| Support knowledge acquisition | Studies find LLM-assistants helpful in learning and knowledge acquisition as a direct [PS6, PS9, PS14, PS28, PS33] and indirect benefit [PS3, PS13]. LLM-assistants are commonly used as an expert consult Other studies envision these assistants taking the role of a tutor in future development workflows [PS21, PS33]. |
| Support code-adjacent tasks | LLM-assistants are found helpful in the ideation process [PS28], requirements specifications [PS28, PS32], documentation [PS1, PS6, PS9, PS12, PS28], and quality assurance [PS9, PS24]. |
| Reduce task initiation overhead | Participants report benefits at the early stages of projects [PS3, PS35] and highlight LLM-assistants' ability to reduce the entry barrier [PS24]. LLM-assistants also support building proof-of-concept applications [PS9]. Developers utilize these tools to generate initial code scaffolding [PS15]. |
| Improve code quality | LLM-assistants have the ability to enhance the quality of code [PS35], which is seen as a key advantage [PS6]. Two studies find improvement in code quality [PS19], with metrics such as cyclomatic complexity, code coverage, technical debt, defect density [PS29], code translation error rate [PS3], code smells [PS29, PS31], defect rate [PS29], and number of defects [PS31]. |
| Support debugging/ troubleshooting | Participants leverage LLM-assistants to help interpret error messages [PS28] and suggest potential fixes [PS20] without the need to consult extensive documentation [PS6]. |

agent (i.e., ChatGPT). The study finds significant productivity gains for both code completion and conversational agents across all five *SPACE* dimensions—satisfaction, performance, activity, communication, and efficiency.

Similarly, [PS11] conducts a controlled experiment showing that using an LLM-assistant plugin for code comprehension leads to statistically significant improvements in task completion rates compared to conventional web search. Similarly, [PS2] introduces a PyCharm plugin designed to reduce reliance on Stack Overflow. However, their study finds no statistically significant differences in task completion time or correctness compared to traditional web browsing, suggesting a difference in the benefits across tools or tasks. Indeed, some studies emphasize that the effectiveness of LLM-assistants over traditional online search is task-dependent. [PS19] conducts a comparative study with 44 participants to evaluate the use of ChatGPT (23 out of 44) versus Stack Overflow (21 out of 44) across programming tasks,

including algorithmic problems, library usage, and debugging. They find that ChatGPT users produce higher-quality code in algorithmic and library usage tasks, whereas the Stack Overflow group perform better in debugging tasks. However, no statistically significant difference was observed between the groups in terms of task completion time.

*6.1.2    Accelerate software development.* Time remains a central consideration in many definitions of developer productivity [31, 32, 59, 60] as time is both a valuable and constrained resource within development workflows. Participants from several empirical studies, particularly those using self-reported methods, suggest that LLM-assistants can accelerate software development [PS9, PS20, PS27, PS35]. Participants often report that LLM-assistants help maintain a state of flow (e.g., "stay in the flow") [PS9, PS23, PS35] and contribute to a perceived increase in productivity [PS35]. Supporting this perception, a qualitative analysis of open-ended feedback on an LLM-based code completion tool [PS12] finds "accelerate coding" to be the second most frequent theme, cited in 14 responses (20%).

Complementing these self-reports, studies using quantitative measures demonstrate that LLM-assistants can reduce task completion time [PS11, PS15, PS25, PS32]. For instance, a case study on the integration of LLMs in the software development lifecycle (SDLC) of a pension plan website [PS32] finds that the required effort decreased from 75 person-days to 22 person-days, representing a productivity gain of 71%. Similarly, statistical significance in time completion has been observed in coding puzzles. This aligns with one of the modes of human–AI interaction described by Barke, James, and Polikarpova [61] as the acceleration mode.

*6.1.3    Automate trivial/ repetitive tasks.* Using LLM-assistants help minimize repetitive coding [PS6, PS9, PS14, PS33] and reduces trivial tasks by generating boilerplate code [PS3, PS9, PS12, PS20, PS31, PS35]. More specifically, some studies report a reduction in keystrokes and typing effort [PS9, PS17]. The broader impact of offloading cognitively repetitive work is further highlighted through a Delphi judgment study conducted with 14 industry professionals to discuss the future of SE in the age of LLM-assistants [PS33]. The study anticipates that LLM-assistants could be used to automate all routine tasks, hence freeing up developer time for more complex tasks.

*6.1.4    Support knowledge acquisition.* The benefits of LLM-assistants extend beyond artifact generation (e.g., source code, test cases). Professional developers increasingly perceive these tools as a valuable aid for learning and knowledge acquisition [PS9, PS28]. LLM-assistants are found to support knowledge acquisition both as a direct [PS6, PS9, PS14, PS28, PS33] and indirect benefit [PS3, PS13].

The authors of [PS3] highlight the concept of knowledge acquisition as an indirect benefit of using LLM-assistants. Although the main goal is to speed up the development process, 69% of participants report that the employed code translation LLM-assistant enhanced their learning (e.g., taught them new aspects of Python) [PS3].

Further evidence is provided by Khojah et al. [PS14], who analyze developer interactions with ChatGPT. Their findings show that expert consultation is the most common use case, accounting for 62% of the analyzed conversations [PS14]. Participants from the same study further highlight such benefits, where 75% of survey respondents report ChatGPT as a helpful learning tool [PS14].

Other studies discuss the potential of using LLM-assistants to teach humans [PS21, PS33]. For instance, an opinion paper about the future of developers in the age of AI highlights the likelihood of LLMs as a tutor as one of the future interaction scenarios [PS21]. Additionally, a Delphi study involving 14 experts in SE finds that enhancing learning and teaching is the most probable future scenario with LLM-assistants [PS33].

*6.1.5    Support code-adjacent tasks.* The benefits of LLM-assistants extend beyond coding tasks to support code-related activities. For instance, studies highlight the use of LLM-assistants for ideation [PS28] by exploring different solutions

options, requirements specifications [PS28, PS32] including functional and non-functional requirements, documentation [PS1, PS6, PS9, PS12, PS28] such as in-code documentation and API documentation, and quality assurance [PS9, PS24].

*6.1.6    Reduce task initiation overhead.*  A common reported benefit across primary studies is the use of LLM-assistants to support task or project initiation [PS3, PS9, PS15, PS24, PS35]. Several studies note that developers rely on these tools as a *starting point* for a project or task [PS3, PS35], effectively lowering the entry barrier [PS24]. These tools help developers build momentum by reducing the time and cognitive effort required during the early stages of a project. For example, developers highlight how LLM-assistants can accelerate the development of proof-of-concept applications by generating multiple candidate implementations for the same task [PS9]. In a qualitative controlled experiment, [PS15] analyzes the interactions of the developers with ChatGPT and finds that 55% of participants (17 out of 31) used the assistant primarily to generate initial code scaffolding. After this initial phase, participants transitioned to more independent workflows by refining and correcting the code themselves and only relying on the LLM-assistant for targeted questions.

*6.1.7    Improve code quality.*  Studies highlight the ability of LLM-assistants to improve code quality [PS3, PS6, PS29, PS31, PS35]. Interview participants in [PS35] report using these tools to rewrite and improve the quality of the code. Similarly, 14% of survey respondents in [PS6] identify improved code quality as a key advantage of LLM-assistants.

These self-reported perceptions are further supported by empirical evidence. [PS29] compares ten projects developed with the support of LLM-assistants and ten developed without such assistance. The authors evaluate six code quality metrics, including cyclomatic complexity, code coverage, code smells, technical debt, and defect density. They report an 18% improvement across all metrics for projects developed with LLM-assistants. Similarly, [PS31] conducts a case study involving five development teams to examine changes in code quality before and after adopting Copilot. Their findings indicate that three of the five teams experience a measurable reduction in code smells, and all five teams show a decrease in the number of software defects following the integration of LLM-assistants (i.e., Copilot), into their development workflows.

Two controlled experiments further highlight these code quality improvements. In the controlled study [PS19], the authors compare the code produced by participants using ChatGPT (i.e., treatment group) with the one produced by those using Stack Overflow (i.e., control group). The results show higher code quality for the ChatGPT group in algorithmic and library usage tasks, although the control group outperformed in debugging tasks. Similarly, [PS3] evaluates code translation quality with and without LLM-assistant (i.e., TransCoder). The authors measure error rates using several translation-related metrics (e.g., Translation Error, Language Error, Spurious Error) and find that the group supported with LLM-assistants exhibits fewer translation errors compared to the control group, resulting in a 51% reduction in error rate.

*6.1.8    Support troubleshooting / debugging.*  Despite the use of earlier generations of bots for automating troubleshooting tasks in software development [PS1], the emergence of LLMs represents the continuation and expansion of this use case. As illustrated in Figure 7 and Table 6, recent studies show that LLM-assistants now play an active role in supporting developers during debugging and troubleshooting [PS28]. These tools help interpret error messages by explaining potential causes and suggesting actionable fixes [PS20]. Additionally, several developers report that LLM-assistants accelerate the debugging process, often eliminating the need to consult extensive documentation [PS6].

Table 7. Summary of reported risks associated with LLM-assistants in developer productivity

| Theme | Summary |
|---|---|
| Limit code quality | Concerns were raised about the quality and accuracy of LLM-generated codes [PS6, PS23]. Vulnerabilities and bugs can be introduced if a developer overestimates the capabilities of the tool [PS5, PS20]. Working with LLM-assistants might not yield better code quality [PS15, PS30]. |
| Fail to meet requirements | LLM-assistants often do not meet functional or non-functional requirements [PS9, PS17]. Developers suggest that not all LLM-assistants' outputs have good accuracy [PS27, PS35]. This is due to the limited controllability of LLMs' output [PS17], as they can be out of context [PS17, PS20] and tend to over-deliver by providing too much information or repetitive code [PS17, PS24]. This calls for iterative refinement to interact more effectively with LLM-assistants [PS27]. |
| Promote over-reliance and cognitive offloading | Several studies raise concerns on the diminishing of critical thinking skills among novices and students [PS8, PS10, PS22, PS37]. In professional settings, instances of automation complacency have also been reported [PS26]. To address these issues, studies recommend promoting cautious and informed LLM-assistants use, emphasizing interactive engagement over passive acceptance [PS15, PS22]. |
| Reduce team collaboration | Relying on LLM-assistants introduces the risk of hindering team collaboration and communication. [PS1, PS14]. This is outlined as one of the limitations of current practices [PS5]. This highlights the need to investigate the impact of LLM-assistance for both human-human and human-agent collaboration and communication. [PS4, PS8]. |
| Disrupt the flow | LLM-assistants can disrupt developers' flow with unwanted suggestions [PS35], interface switching [PS24], verbose answers [PS24], and inadequate speed of code suggestions [PS24]. The simultaneous use of multiple code completion tools can disrupt developer flow, particularly when competing suggestions are presented [PS12]. |

## 6.2 Risks

Table 7 shows several risk themes identified across the primary studies that may affect developer productivity. In this section, we describe the five risks categories: limit code quality, fail to meet requirements, promote over-reliance and cognitive offloading, reduce team collaboration, and disrupt the flow.

*6.2.1 Limit code quality.* Code quality is the one theme that has been reported as both a benefit and a risk (see Figure 7, Table 6 and 7). Multiple studies raise concerns about the quality of generated code. For instance, [PS6] reports that 13% of survey respondents raise concerns about the quality and accuracy of LLM-generated code. Similar concerns are raised by 29% of survey participants when using Copilot [PS23]. Code quality issues arise when developers overestimate the capabilities of such tools, which can introduce vulnerabilities and bugs [PS5]. For example, LLM-assistants often struggle with optimization and refactoring tasks, especially when lacking semantic context [PS5]. Moreover, in [PS20], 38% of the surveyed participants point to erroneous code as one of the limitations of ChatGPT, while 29% report limitations due to inefficient code. User experiments conducted by the authors of [PS15] reveal no significant improvements in code quality or correctness with LLM-assistants, with a slightly worse quality average for the ChatGPT group.

When analyzing the relationship between reported productivity and code quality gains in the context of LLM-assistant, a large industry study involving 70 large global companies [PS30] reports a negative correlation between the two. These findings highlight that increased productivity through the support of LLM-assistants does not necessarily lead to improvements in code quality.

*6.2.2   Fail to meet requirements.*  Developers acknowledge that not all the suggestions of LLM-assistants are accurate [PS35]. Many survey participants in [PS9, PS17] mention that LLM-assistants often fail to meet both functional and non-functional requirements. Some developers perceive them as difficult to control [PS17], noting instances where responses are out of context [PS17, PS20] or tend to over-deliver by providing too much information or repetitive code [PS17, PS24]. For example, 50% of participants report missing or misunderstanding the requirement context as the two main issues encountered with ChatGPT 3.5 [PS20]. Generic or inaccurate code suggestions often require additional effort to modify and refine, leading to an iterative process of prompt refinement and learning how to interact effectively with LLM-assistants [PS27].

*6.2.3   Promote over-reliance and cognitive offloading.*  A heavy reliance and excessive trust in LLM-assistants raises concerns about the erosion of critical thinking skills, especially for novice developers and students [PS8, PS10, PS22, PS37]. For instance, authors of [PS37] develop an AI tutor that limits direct interaction with ChatGPT through predefined prompts, aiming to promote critical thinking and reduce dependence on LLM-assistants for every minor challenge. However, students still expressed concerns post-experiment, noting that reliance on the LLM tutor might hinder their learning progress [PS37]. The authors acknowledge this issue and highlight the need for further refinement of the tool.

Over-reliance and automation complacency have also been documented among professional software engineers. In a study involving a programming exam with Google Bard, [PS26] observes all three characteristics of automation complacency, as described by Parasuraman and Manzey [62]: human monitoring of an automated system, infrequent monitoring, and degraded performance. These findings highlight the need to promote responsible LLM usage. Several studies advocate for more interactive and reflective engagement with LLM outputs [PS15], cautioning against blind trust in automated responses [PS22] and encouraging users to understand the tools' capabilities and limitations.

*6.2.4   Reduce team collaboration.*  Relying on LLM-assistants can negatively impact productivity by reducing team collaboration and communication [PS1, PS14]. For instance, a field study [PS14] observes that excessive use of LLM-assistants may lead developers to favor consulting a chatbot over a colleague. In fact, the overconfidence of LLM-assistants' responses can create the impression that team discussions are unnecessary, reducing opportunities for communicative learning and discovery [PS14]. The intrinsic nature of LLM-assistants also plays a role in reduced team collaborations. For instance, current conversational LLM-assistants offer limited support for team collaboration, as they primarily support one-on-one interactions and are not well-suited to facilitating effective team coordination [PS5].

These findings highlight the need for future studies to further investigate how LLM-assistants affect human-human collaboration and how they can be designed to foster team collaboration [PS4, PS8].

*6.2.5   Disrupt the flow.*  Studies find that LLM-assistants can disrupt developer flow [PS12, PS35]. Issues that impact developer state of flow have been attributed to various kinds of interruptions, including unwanted LLM suggestions [PS35], interface switching, and verbose answers [PS24]. For example, in a laboratory experiment [PS24], some professional developers find Copilot distracting, as the speed of code suggestions does not allow sufficient time for code understanding. Distraction also occurs when LLM-assistants work in tandem and compete to display suggestions [PS12].

Authors of [PS23] investigate developers' interaction with code recommendation systems and their impact on flow by modeling user behavior while using tools such as Copilot. Findings show that developers spend an average of 51.5% of their coding session time in LLM interaction states, such as verifying suggestions, prompt crafting, and deferring thought. These findings highlight the temporal and cognitive costs that these tools may introduce.

---

**RQ 2 - Summary**

Studies report mixed findings on the use of LLM-assistants, revealing both notable benefits and critical risks. The most frequently reported benefits include minimizing the need for online code search, accelerating development, and automating trivial or repetitive tasks. At the same time, primary studies identify several risks, such as promoting over-reliance and cognitive offloading, limiting code quality, and reducing team collaboration. Code quality emerges as a particularly contested area, with evidence pointing to both improvements and degradations depending on context. These discrepancies underscore the need for further investigation and the development of strategies to ensure that code quality is maintained when integrating LLM-assistants into software development workflows.

---

## 7 RQ3: WHICH DIMENSIONS OF DEVELOPER PRODUCTIVITY ARE INVESTIGATED, AND HOW DO THESE DIMENSIONS MAP ONTO THE SPACE FRAMEWORK?

The diverse methodologies and wide variation in reported findings across primary studies (see Section 5 and Section 6) highlight the need for a structured lens to interpret how productivity is conceptualized in the literature, especially in the context of rapidly evolving development with LLM-assistants. We leverage the *SPACE* framework [19], a widely used framework developed by Microsoft researchers, as an organizing model to understand productivity from a practical perspective. We choose the *SPACE* framework because it defines productivity as a multi-dimensional construct that covers a diverse range of study instrumentation from measurement metrics (e.g., code quality) to perceptions (e.g., satisfaction). The *SPACE* framework also reflects a modern development workflow, aligned with the complexities of LLM-assisted development, which emphasizes dimensions like collaboration and trust. Its extensible and adaptable nature allows for contextual adaptation, as it does not prescribe fixed metrics but offers a conceptual structure that can be tailored to specific empirical contexts. We map each primary study to the five dimensions of the *SPACE* framework: Satisfaction, Performance, Activity, Communication & collaboration, and Efficiency (see Figure 8).

To synthesize and compare findings from our primary studies, we leverage thematic analysis, employing an adapted version of the *SPACE* framework. We use the five dimensions of the *SPACE* framework. To provide more granularity, we further refine these dimensions by including sub-dimensions. Specifically, we adapt some sub-dimensions from relevant related work [63]. When a specific concept was not captured by an existing sub-dimension, we ensure comprehensive coverage by including emerging sub-dimensions derived from our data. A structured taxonomy for our analysis is provided in Table 8. Figure 8 illustrates the distribution of the sub-dimensions across each dimension of the *SPACE* framework.

Our findings reveal that the majority of studies (92%, 34 out of 37) adopt a multidimensional perspective on productivity, with only three studies taking a uni-dimensional perspective. This indicates a shift away from singular-dimension perspectives toward a more holistic understanding of the complex nature of productivity in software engineering (See figure 9). However, only 14% of the studies (5 out of 37) examine more than three *SPACE* dimensions, with just four addressing all five dimensions and only one study covering four. This highlights the need for future work to capture the full breadth of how LLM-assistants impact productivity. We find that the most co-occurring combinations

Table 8. Mapping of primary studies to SPACE dimensions and derived sub-dimensions. With sources for each sub-dimension are cited from prior work.

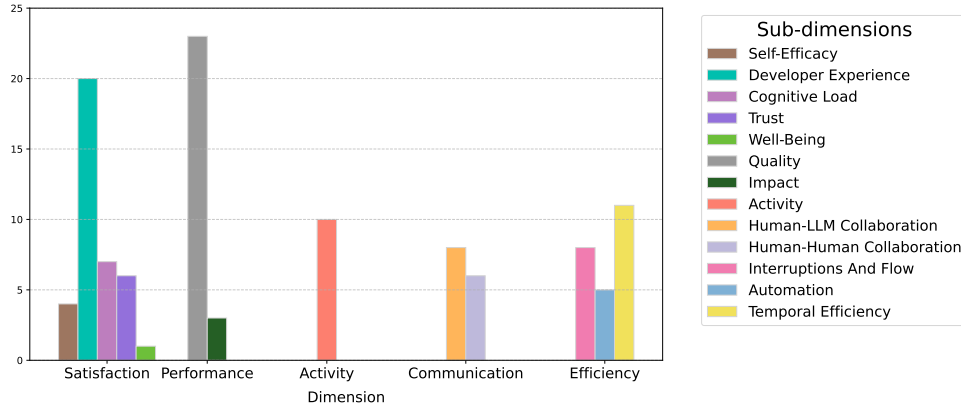| Dimension | Sub-dimensions | Primary Studies | % |
|---|---|---|---|
| Satisfaction | Developer experience [26] | [PS2, PS6, PS9, PS11, PS12, PS14, PS16, PS17, PS18, PS19, PS20, PS22, PS23, PS24, PS25, PS27, PS28, PS29, PS31, PS37] | 78% |
| | Self-efficacy [19, 63] | [PS4, PS10, PS35, PS36] | |
| | Trust [63] | [PS4, PS10, PS14, PS20, PS21, PS26] | |
| | Well-being [19] | [PS5] | |
| | Cognitive load | [PS3, PS10, PS11, PS15, PS25, PS26, PS36] | |
| Performance | Quality [19, 63] | [PS1, PS2, PS3, PS4, PS5, PS7, PS10, PS12, PS15, PS16, PS19, PS20, PS22, PS23, PS24, PS25, PS26, PS29, PS30, PS31, PS32, PS34, PS36] | 65% |
| | Impact [19] | [PS30, PS33, PS34] | |
| Activity | | [PS7, PS10, PS12, PS13, PS16, PS18, PS22, PS23, PS24, PS31] | 27% |
| Communication | Human-LLM collaboration | [PS8, PS11, PS16, PS21, PS23, PS24, PS28, PS37] | 35% |
| | Human-human collaboration | [PS1, PS4, PS5, PS8, PS14, PS31] | |
| Efficiency | Temporal efficiency [19] | [PS2, PS11, PS13, PS15, PS19, PS24, PS26, PS29, PS30, PS31, PS32] | 59% |
| | Interruptions and flow [19, 63] | [PS4, PS9, PS16, PS21, PS23, PS27, PS33, PS35] | |
| | Automation | [PS1, PS6, PS9, PS18, PS23] | |



Fig. 8. Distribution of sub-dimensions across each dimension of the *SPACE* framework.

involve Satisfaction, Performance, and Efficiency. Most studies covered pairs of dimensions (16 out of 37), and the most frequent pairs are Satisfaction-Efficiency (4 out of 37) and Satisfaction-Performance (4 out of 37).

Table 8 and Figure 9 show that **Satisfaction** is the most studied dimension, addressed by 78% (29 out of 37) of primary studies. This dimension captures developers' feelings about their work with LLM-assistants, which is mainly captured through self-reported instruments. Our analysis of this dimension reveals five fine-grained sub-dimensions (i.e., *developer experience*, *self-efficacy*, *trust*, *well-being*, and *cognitive load*). Most studies within the satisfaction dimension focus on the concept of *developer experience*, which encompasses developers' perceptions, feelings, and values regarding
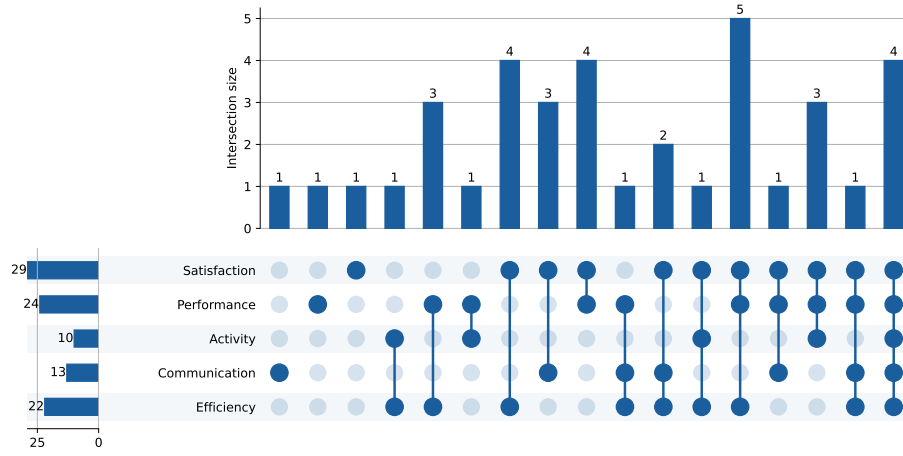
Fig. 9. The distribution of investigated SPACE dimensions and their overlap.

their interactions with LLM-assistants [26], as well as the perceived importance of LLM-assistants and their ease of use (e.g., developers' feedback or the Technology Acceptance Model (TAM)). *Cognitive load* is the second most commonly investigated satisfaction sub-dimension, assessed using instruments such as NASA-TLX or custom surveys. *Self-efficacy*, the belief in one's ability to complete tasks, is explored using both validated and custom-designed tools. Finally, *Well-being* is notably underexplored, appearing in only one non-empirical study [PS5], which highlights how developers' well-being and mental health are often overlooked in the context of software engineering.

**Performance** is the second most studied dimension covered by 68% (24 out of 37) of primary studies, and often together with satisfaction. Performance mainly concerns the final outcomes of software development activities. The majority of studies investigate the *quality* sub-dimension of performance using instruments such as passing unit tests, functional correctness, and code smells (see Table 9). The *impact* sub-dimension is addressed in only three studies and mainly investigates how LLM-assistants impact final product outcomes [19]. This sub-dimension focuses on business-related metrics, including cost savings, product quality improvements, and delivery speed (see section 5.3.4).

**Efficiency** is examined by 59% of studies (22 out of 37) and often in tandem with satisfaction. This dimension reflects the capacity to complete tasks efficiently with minimal interruptions or time delays, as it aims to minimize unnecessary delays and optimize the flow of task handoffs [19]. We highlight three sub-dimensions explored by primary studies (i.e., *temporal efficiency, automation, interruptions and flow*). Efficiency is often investigated from the temporal perspective, measured using task completion metrics or via developer perceptions. *Automation* is another important angle of efficiency, with studies reporting how LLM-assistants are used to offload repetitive tasks such as writing boilerplate code [PS6, PS9]. Finally, studies examine efficiency in terms of *interruptions and flow*, with some highlighting reduced cognitive interruptions and others noting new forms of distraction introduced by the LLM-assistant.

The **Activity** dimension is the least investigated dimension across primary studies (24%, 10 out of 37). Activity is often paired with efficiency and performance and focuses on counts and frequency measures while performing a given task Forsgren et al. [19]. Studies included in our SLR often measure activity as the count of actions or tasks

developers take during their interactions with LLMs (e.g., acceptance rate, number of tasks completed)[PS12, PS13, PS16, PS24]. For example, Ziegler et al. [PS16] measures the dimension of activity at a finer granularity, regarding the count of actions developers take during their interactions with Copilot (e.g., acceptance rate, suggestions shown rate, completions changed or unchanged).

**Communication** is one of the least explored dimensions (35%, 13 out of 27). Communication focuses on how developers and teams communicate and share knowledge [19]. Most empirical studies investigate communication in terms of *human-LLM collaboration*, which includes analysis of interaction patterns between participants and LLM-assistants. In contrast, very few studies (6 out of 13) examine *human-human collaboration* with LLM in the loop. This highlights a gap in our understanding of how LLM-assistants influence team communication or coordination. Given that emerging concerns regarding reduced team collaboration due to over-reliance on LLM-assistants (see Section 6.2.4), future studies should incorporate more investigation on team dynamics to better understand their impact in the LLM-assisted workflows.

Table 9. Quality metrics by study.

| Metric | Primary Studies |
| --- | --- |
| Passing Unit Tests | [PS3, PS15, PS19, PS20, PS36] |
| Functional Correctness and Accuracy | [PS4, PS7, PS19, PS32] |
| Code Smells | [PS10, PS29, PS31] |
| BLEU Score | [PS7, PS12] |
| Halstead Complexity Measures | [PS24, PS32] |
| Cyclomatic Complexity | [PS2, PS29] |
| Translation Error Rate | [PS3] |
| Maintainability Index | [PS24] |
| Cognitive Complexity | [PS31] |
| Defect Density | [PS29] |
| Defect Rate | [PS31] |
| Technical Debt | [PS29] |
| Code Coverage | [PS29] |

**RQ3 - Summary**

Our analysis, framed by the *SPACE* framework, reveals that the majority of studies (92%) have a multidimensional view of productivity in SE, mostly combining two or more dimensions of the *SPACE* framework. However, only 14% examine more than three *SPACE* dimensions. Satisfaction (78%), and Performance (65%) are the most investigated productivity dimensions. In contrast, Activity (24%) and Communication (35%) are the least explored productivity dimensions.

## 8 DISCUSSION

In this section, we discuss how LLM-assistants impact software development through the theoretical framework of McLuhan's Tetrad [64]. We analyze how LLM-assistants enhance existing tools, make obsolete, retrieve, and reverse some practices. We also present practical and research implications, with an emphasis on trust, workflow changes, productivity challenges, and recommendations for practitioners and researchers.
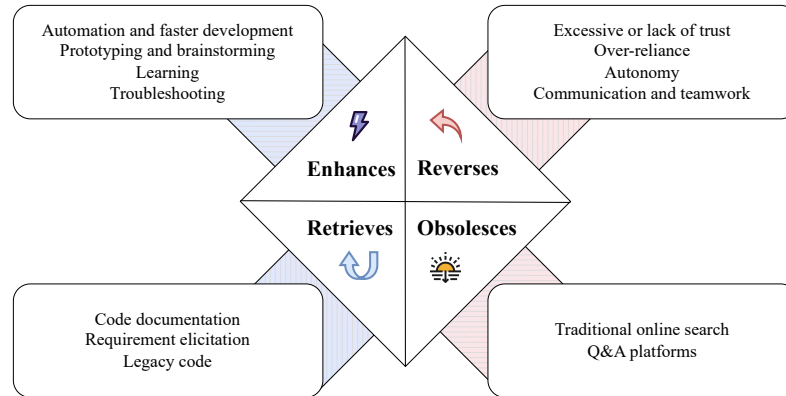
Fig. 10. McLuhan's Tetrad diagram illustrates the implications of LLM assistants on the productivity of software developers. The diagram captures four dimensions. *Enhancement*: how LLM-assistants amplify development speed; *Obsolescence*: which traditional practices are being displaced; *Retrieval*: which previously diminished practices are being revived; and *Reversal*: what adverse effects may emerge when LLM-assistants are pushed to the extreme.

## 8.1 LLM-assistants in software development: McLuhan Tetrad

Before the emergence of LLM-assistants, researchers had long emphasized the challenges of conceptualizing and measuring productivity in SE [18]. Traditional productivity metrics have been criticized for their negative consequences, including distorted incentives, decreased employee trust, and ineffective management decisions [65]. The integration of LLM-assistants is also extended to the entire software development life-cycle [PS14, PS31, PS32], introducing new dynamics that further complicate these measurements, simultaneously amplifying certain activities and diminishing the relevance of others. As discussed in Section 6.1, while LLM-assistants bring several productivity-enhancing benefits, such as improved code quality and accelerated development, they also raise challenges related to over-reliance and reduced team collaboration as discussed in Section 6.2.

Providing a clear understanding of how LLM-assistants affect the software development process is therefore essential to harness their full potential. Moreover, their utilization in the day-to-day development practices prompts broader questions regarding the future skill sets required of professional developers and the evolving nature of SE [PS8]. To interpret these shifts more holistically, we draw on McLuhan's tetrad [64], a framework that offers a structured approach to examine the disruptive effects of emerging media on society. The tetrad poses four interrelated questions: 1) *What does the medium enhance?*, 2) *What does it render obsolete?*, 3) *What does it retrieve?*, and 4) *What does it reverse when pushed to the extreme?* Figure 10 illustrates how this framework is applied to analyze the impact of LLM-assistants on software development.

*8.1.1 Enhance.* According to McLuhan's Tetrad, the enhanced dimension refers to what technology intensifies in existing practices. In the realm of SE, LLM-assistants have the potential to enhance several aspects of the development process. Several studies highlight the ability of LLM-assistants to streamline labor-intensive tasks [PS14], thereby improving development speed. For example, LLM-assistants can solve boilerplate tasks and implement basic algorithms [PS17]. In early development phases, LLM-assistants can help to build proof-of-concept applications [PS9] or brainstorm

architectural ideas [PS14]. One participant describes such a tool as a "brainstorming engine" Beyond task automation, LLM-assistants also enhance learning and knowledge acquisition. In a SE educational context, these tools can provide timely feedback and guidance for students [PS37]. For professional developers, LLM-assistants offer support for exploring unfamiliar technologies, troubleshooting, and debugging, ultimately contributing to skill development and continuous learning [PS20].

*8.1.2 Reverse.* The reverse dimension in McLuhan's Tetrad refers to the unintended or counterproductive consequences that may emerge when technology is pushed to its limits. A primary concern is the polarization of trust. On one end, some developers exhibit excessive trust, leading to over-reliance and automation complacency [PS22, PS36], which can lead to compromised code quality. This over-reliance on automation may increase the risk of reduced communication and awareness between developers [PS8]. On the other end, the under or complete lack of trust in LLM-assistants may result in frustration and disengagement [PS3, PS10], limiting the tools' effectiveness and reducing the likelihood of adoption. Low trust in LLM-assistants may be associated with a desire to preserve the sense of control and autonomy during development tasks [PS8, PS21].

*8.1.3 Obsolesce.* The obsolesce dimension of McLuhan's Tetrad refers to the tools, practices, or roles that are diminished or displaced as a result of adopting new technologies. In the context of LLM-assistants, one aspect at risk of becoming obsolete is traditional online search techniques for recalling syntax [PS9, PS17], such as reading documentation or looking up new libraries online [PS17]. Developers may prefer using conversational agents over searching for solutions or consulting Q&A platforms such as Stack Overflow [PS6, PS8, PS9, PS17, PS35]. Furthermore, manual documentation, often perceived as a labor-intensive task, is increasingly being augmented or replaced by LLM-generated summaries [PS12].

*8.1.4 Retrieve.* In McLuhan's Tetrad, the retrieve dimension reflects the resurgence or reintegration of practices that had diminished in relevance prior to the adoption of new technology. In the context of LLM-assistants, several studies capture the reemergence of previously neglected practices or aspects of software development. One such practice is code documentation, which is now being brought back to focus in development workflows with the help of LLM-assistants [PS1, PS6, PS9, PS12] after being de-prioritized in fast-paced development environments. Another aspect is requirements engineering, where tasks such as requirement elicitation, traditionally seen as time-consuming and resource-intensive due to the need for frequent client communication [66], are gaining momentum through LLM-assistants [PS5, PS6, PS32]. Additionally, LLM-assistants may help revive attention to legacy systems, a domain often overlooked due to the high cost and complexity of rewriting code in outdated or niche languages such as COBOL or Uniface. While some studies highlight industrial use cases where LLM-assistants are involved in modernization efforts [PS27, PS33], others report that these tools currently provide limited support for legacy systems, especially those built on less common platforms [PS17]. This suggests a promising but underexplored opportunity for future development and tool improvement.

## 8.2 Implications and recommendations for software developers

The integration of LLM-assistants into the software development workflow has profound implications, transforming both the individual developer's role and the team's dynamics. While these tools offer benefits, developers must be mindful of the potential risks and adapt their practices to maximize productivity and maintain high-quality code.

**Cultivating and managing trust.** Establishing an appropriate level of trust between developers and LLM-assistants presents its challenges [PS12]. Trust can become fragile when these models hallucinate or produce incorrect outputs,

such as suggesting non-existent APIs or generating syntactically incorrect code Murali et al. [PS12]. These issues stem from the traditional trade-off between precision and recall in machine learning systems, where models tuned for high recall may increase coverage but at the cost of higher volume of incorrect suggestions Murali et al. [PS12]. This creates a trust barrier between developers and LLM-assistants.

Additionally, LLMs often suffer from a lack of transparency, as they do not provide sources or references for their outputs [PS14] which deepens the trust gap. Implementing AI solutions with transparency mechanisms [PS6] to elucidate their decision-making processes empowers developers to trust and better understand the LLM-assistants' suggestions and decisions. These implications may extend to a much broader scope. Through a questionnaire conducted by [PS20], the authors find a moderate positive correlation between self-reported productivity and trust levels, as participants who report increased productivity with LLM-assistants also exhibit slightly higher levels of trust. However, trust dynamics are not uniform across user groups. For instance, the study in [PS26] finds that novice developers tend to demonstrate automation complacency, often accepting LLM-assistants suggestions uncritically and with minimal validation. While this may temporarily boost productivity, it also raises concerns about long-term skill development and critical thinking.

> *Recommendation 1.* Developers should cultivate calibrated trust by understanding the capabilities and limitations of LLM-assistants. Fostering informed, critical trust is essential to maximizing benefits without compromising code quality, autonomy, or learning. Actively verifying outputs and ensuring human oversight can help maintain a productive level of trust.

**Redefining the developer's role from coder to reviewer.** LLM-assistants often lack awareness of the broader context and intricacies of a complex software project [PS9, PS17]. Developers often mitigate this issue by breaking down the problem [PS9, PS17, PS20] and providing a clear explanation to the LLM-assistants [PS9, PS17]. Consequently, developers increasingly spend time verifying, editing or refining LLM-assistants generated suggestions rather than writing code. One study by Mozannar et al. [PS23] reports that participants spend over 50% of their time in evaluation activities, such as crafting prompts, reviewing suggestions, and editing completion. Similarly, Weisz et al. [PS3] finds that using LLM-assistants for code translation transforms the developer's responsibility into one of reviewing rather than writing code. Time saved in code generation can be lost in the evaluation and refinement phases, especially for complex tasks. This can lead to diminishing returns on productivity if not managed effectively. [PS30]. These findings align with related work regarding "ironies of automation" [67], where the introduction of automation often shifts users' roles from production to evaluation.

> *Recommendation 2.* The role of developers is evolving from coder to reviewer when using LLM-assistants. To maximize productivity and code quality, developers must shift their mindset and reallocate time and cognitive resources to tasks like prompt engineering, iterative evaluation, and refining LLM-generated outputs. This includes guiding LLMs with precise context, treating suggestions as drafts requiring human validation, and decomposing complex tasks into smaller, well-scoped sub-problems for more coherent responses [PS9, PS17, PS20].

**Adapting the development workflow.** The integration of LLM-assistants may disrupt established software development workflow, which can affect team collaboration and individual tasks. At the team level, we describe in section 6.2.4 how LLM-assistants may reduce collaboration among software teams. It remains unclear to what extent the

adoption of LLM-assistants leads to a decline in developer knowledge sharing and pair programming. At the individual level, LLM-assistants can disrupt developers' flow. As described in section 6.2.5, unwanted or irrelevant LLM-assistants' suggestions can disrupt the developer's flow. For instance, studies report that interruptions from tools like Copilot, especially when suggestions are too frequent or conflict with other tools, can disrupt focus and reduce productivity [PS12, PS24, PS35].

> *Recommendation 3.* Developers must adapt their practices and team dynamics to ensure LLM-assistants enhance, rather than hinder, the development workflow. Individually, developers should customize tool settings to control suggestion frequency and limit interruptions. At the team level, preserving collaborative practices like pair programming, code reviews, and architectural discussions is crucial to maintaining communication and knowledge sharing.

### 8.3 Implications and recommendations for researchers

**Establishing shared practices and cumulative insight** The current body of research is primarily exploratory and formative (64%). Laboratory experiments represent the most common research strategy (41%), focusing on controlled tasks to isolate effects. This is critical for early insights. However, the methodological diversity can challenge cross-study comparison and synthesis. For instance, LLM-assistant's effect on code quality is inconsistent, reported as both a benefit and a risk depending on study context and metrics used. Similarly, cognitive load findings are mixed: some studies report reduced mental effort, while others find no effect or increased frustration. These inconsistencies highlight the complexity of the research space and the need for more context-aware, longitudinal studies that account for individual, task, and organizational differences.

> *Recommendation 1.* Researchers should adopt shared evaluation frameworks, validated instruments that allow comparability. Future work should incorporate field studies and team-based studies and triangulate quantitative and qualitative data to better capture context aware trade-offs.

**On the dimensions of productivity** Existing work agrees that developer productivity is a multi-dimensional and context-sensitive construct [25]. This complexity is magnified in the context of LLM-assisted development. Our findings provide evidence that the vast majority of studies (92%) examine at least two *SPACE* dimensions. This indicates a positive move from the research community toward a multidimensional perspective. However, our findings also highlight that only 14% of studies extend beyond three dimensions, indicating that room for improvement remain in terms of providing a more complete evaluation of the productivity concept.

Satisfaction, Performance, and Efficiency are the most frequently investigated dimensions. Communication, however, remains underexplored with the human-human collaboration being the least investigated communication sub-dimension. This is inline with our findings on LLM-assistants risks on collaboration and further highlights the need to investigate the impact of LLM-assistance for both human-human and human-agent collaboration and communication.

*Recommendation 2.* Researchers should continue advancing multidimensional evaluations of developer productivity by systematically addressing underexplored *SPACE* dimensions, particularly Communication and Collaboration. While Satisfaction, Performance, and Efficiency are frequently assessed, the human-human and human-agent interaction aspects remain limited. Future studies should incorporate richer measures of team dynamics and communication patterns to better understand how LLM-assistants affect collaborative workflows.

## 9   THREATS TO VALIDITY

**Study selection bias.** A key threat lies in the potential omission of relevant studies due to the inclusion and exclusion criteria defined in our protocol (see Section 3.1.1). Specifically, we exclude shorter papers, non-peer-reviewed work, and publications outside academic journals and conference proceedings. To mitigate this, all authors collaboratively agree on these criteria to ensure methodological rigor and quality control to enhance the reliability of the synthesized findings. Another threat relates to the challenge of identifying the human-centered studies within the large body of LLM research. Our initial search strings include terms such as "performance" or "efficiency", which yielded results focused on the technical applications of LLMs rather than their impact on developer productivity. This issue stems from the broad scope of LLM4SE research [1], where such keywords frequently describe model behavior or algorithmic improvements rather than developer-centric outcomes. To mitigate this threat, three authors jointly review the selected control papers to ensure alignment with the inclusion criteria. We iteratively validate our search strings using control articles inspired by Zhang, Babar, and Tell [45] (see section 3.1.2), ensuring all control papers are correctly retrieved. Furthermore, we employ backward and forward snowballing to further enhance the coverage of relevant articles and reduce the risk of missing key studies.

**SLR bias and repeatability.** The open-ended nature of our research questions can introduce the threat of taking subjective decisions in the interpretation and selection of studies. To mitigate this, we adopt a multi-step validation process involving all co-authors. While the initial screening and data extraction were conducted by one author, the remaining authors were actively involved in selection validation and protocol design. The teams held regular weekly meetings for a period of 9 months to discuss the inclusion decisions and refine the selection and data extraction process. Additionally, we employ a conservative screening strategy during the full-text review, whereby any study with uncertain eligibility was retained for further assessment and independently reviewed by two senior co-authors to minimize the risk of excluding relevant studies.

**Classification rigor.** A potential threat lies in the mapping of study findings to the *SPACE* framework (see Section 7). Since the *SPACE* framework was not originally developed to represent productivity in human-LLM collaboration settings, assigning findings to specific sub-dimensions required interpretive decisions. This may have introduced subjective bias during data coding. To mitigate this, we adapted definitions from prior literature [19, 63] to better fit the context of our review, and discussed coding decisions collaboratively in team meetings.

## 10   CONCLUSION

In this paper, we investigate LLM-assistants' impact on developer productivity. To achieve this, we systematically identify and analyze 37 peer-reviewed studies from their methodological strategies, evaluation practices, and the productivity dimensions these studies focus on. We synthesize reported benefits and risks, and apply established conceptual frameworks to map our findings and contextualize their broader implications.

Our analysis reveals a range of reported benefits, including reduced task initiation overhead, accelerated development, and support for code-adjacent tasks. At the same time, studies identify several risks, such as over-reliance on LLM-assistants especially affecting novice programmers, disruptions to developer flow, and reduced team communication or collaboration. Code quality, in particular, has a mixed outcome, with studies reporting both improvements and degradations depending on context, task design, and evaluation criteria. Our findings show that most studies (92%) consider multiple productivity dimensions. However, relatively few studies extend beyond three, and dimensions like communication and, more specifically, human-human collaboration remain underexplored.

Looking ahead, there is value in expanding the evidence base through team-based studies that capture the dynamic and socio-technical nature of software development. As LLM-assistants become more deeply embedded in everyday workflows, future research will play a critical role in understanding the multidimensional impact of LLM-assistants on developer productivity. To facilitate transparency and future work, we provide a publicly available replication package [20].

## REFERENCES

[1] Xinyi Hou et al. "Large language models for software engineering: A systematic literature review". In: *ACM Transactions on Software Engineering and Methodology* 33.8 (2024), pp. 1–79.

[2] OpenAI. *GPT-4 Technical Report*. Accessed: 2025-06-12. 2023. URL: https://openai.com/research/gpt-4.

[3] GitHub. *GitHub Copilot*. https://github.com/features/copilot. Accessed: 2025-06-12. 2021.

[4] Alessio Buscemi. "A comparative study of code generation using chatgpt 3.5 across 10 programming languages". In: *arXiv preprint arXiv:2308.04477* (2023).

[5] Xiaodong Gu et al. "On the effectiveness of large language models in domain-specific code generation". In: *ACM Transactions on Software Engineering and Methodology* 34.3 (2025), pp. 1–22.

[6] Seohyun Kim et al. "Code prediction by feeding trees to transformers". In: *2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE)*. IEEE. 2021, pp. 150–162.

[7] Yuwei Zhang et al. "PATCH: Empowering Large Language Model with Programmer-Intent Guidance and Collaborative-Behavior Simulation for Automatic Bug Fixing". In: *ACM Transactions on Software Engineering and Methodology* (2025).

[8] Guoyang Weng and Artur Andrzejak. "Automatic bug fixing via deliberate problem solving with large language models". In: *2023 IEEE 34th International Symposium on Software Reliability Engineering Workshops (ISSREW)*. IEEE. 2023, pp. 34–36.

[9] Runchu Tian et al. "Debugbench: Evaluating debugging capability of large language models". In: *arXiv preprint arXiv:2401.04621* (2024).

[10] Shubhang Shekhar Dvivedi et al. "A comparative analysis of large language models for code documentation generation". In: *Proceedings of the 1st ACM International Conference on AI-Powered Software*. 2024, pp. 65–73.

[11] Maram Assi, Safwat Hassan, and Ying Zou. "LLM-Cure: LLM-based Competitor User Review Analysis for Feature Enhancement". In: *ACM Trans. Softw. Eng. Methodol.* (June 2025). ISSN: 1049-331X. DOI: 10.1145/3744644. URL: https://doi.org/10.1145/3744644.

[12] Junjie Wang et al. "Software testing with large language models: Survey, landscape, and vision". In: *IEEE Transactions on Software Engineering* (2024).

[13]  Angela Fan et al. "Large language models for software engineering: Survey and open problems". In: *2023 IEEE/ACM International Conference on Software Engineering: Future of Software Engineering (ICSE-FoSE)*. IEEE. 2023, pp. 31–53.

[14]  Shantanu Mandal et al. "Large language models based automatic synthesis of software specifications". In: *arXiv preprint arXiv:2304.09181* (2023).

[15]  Jules White et al. "Chatgpt prompt patterns for improving code quality, refactoring, requirements elicitation, and software design". In: *Generative ai for effective software development*. Springer, 2024, pp. 71–108.

[16]  Christian Bird et al. "Taking Flight with Copilot: Early insights and opportunities of AI-powered pair-programming tools". In: *Queue* 20.6 (2022), pp. 35–57.

[17]  Emerson Murphy-Hill et al. "What predicts software developers' productivity?" In: *IEEE Transactions on Software Engineering* 47.3 (2019), pp. 582–594.

[18]  Moritz Beller et al. "Mind the gap: on the relationship between automatically measured and self-reported productivity". In: *IEEE Software* 38.5 (2020), pp. 24–31.

[19]  N Forsgren et al. *The SPACE of developer productivity: There's more to it than you think. acmqueue 19 (1), 20–48.* 2021.

[20]  Amr Mohamed, Maram Assi, and Mariam Guizani. *Exploring the Impact of AI and LLMs on Software Developers' Productivity: A Systematic Review and Mapping of the Literature.* Supplemental Material. 2025. DOI: 10.5281/ zenodo.15788502. URL: https://zenodo.org/records/15788502.

[21]  Frederick P Brooks Jr. *The mythical man-month: essays on software engineering*. Pearson Education, 1995.

[22]  Audris Mockus, Roy T Fielding, and James D Herbsleb. "Two case studies of open source software development: Apache and Mozilla". In: *ACM Transactions on Software Engineering and Methodology (TOSEM)* 11.3 (2002), pp. 309–346.

[23]  Stefan Wagner and Melanie Ruhe. "A systematic review of productivity factors in software development". In: *arXiv preprint arXiv:1801.06475* (2018).

[24]  Daniel Graziotin et al. "What happens when software developers are (un) happy". In: *Journal of Systems and Software* 140 (2018), pp. 32–47.

[25]  Kai Petersen. "Measuring and predicting software productivity: A systematic map and review". In: *Information and Software Technology* 53.4 (2011), pp. 317–343.

[26]  Abdul Razzaq et al. "A Systematic Literature Review on the Influence of Enhanced Developer Experience on Developers' Productivity: Factors, Practices, and Recommendations". In: *ACM Computing Surveys* 57.1 (2024), pp. 1–46.

[27]  André N Meyer et al. "Software developers' perceptions of productivity". In: *Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering*. 2014, pp. 19–29.

[28]  André N Meyer et al. "The work life of developers: Activities, switches and perceived productivity". In: *IEEE Transactions on Software Engineering* 43.12 (2017), pp. 1178–1193.

[29]  André N Meyer, Thomas Zimmermann, and Thomas Fritz. "Characterizing software developers by perceptions of productivity". In: *2017 ACM/IEEE International Symposium on Empirical Software Engineering and Measurement (ESEM)*. IEEE. 2017, pp. 105–110.

[30]  Tony Savor et al. "Continuous deployment at Facebook and OANDA". In: *Proceedings of the 38th International Conference on software engineering companion*. 2016, pp. 21–30.

[31]    Minghui Zhou and Audris Mockus. "Developer fluency: Achieving true mastery in software projects". In: *Proceedings of the eighteenth ACM SIGSOFT international symposium on Foundations of software engineering.* 2010, pp. 137–146.

[32]    Capers Jones. "Software metrics: good, bad and missing". In: *Computer* 27.9 (1994), pp. 98–100.

[33]    Moritz Beller et al. "What's DAT? Three Case Studies of Measuring Software Development Productivity at Meta With Diff Authoring Time". In: *arXiv preprint arXiv:2503.10977* (2025).

[34]    Abi Noda et al. "DevEX: What actually drives productivity?" In: *Communications of the ACM* 66.11 (2023), pp. 44–49.

[35]    Fabian Fagerholm and Jürgen Münch. "Developer experience: Concept and definition". In: *2012 International Conference on Software and System Process (ICSSP).* 2012, pp. 73–77. DOI: 10.1109/ICSSP.2012.6225984.

[36]    Darja Smite et al. "Changes in perceived productivity of software engineers during COVID-19 pandemic: The voice of evidence". In: *Journal of Systems and Software* 186 (2022), p. 111197.

[37]    Lan Cheng et al. "What improves developer productivity at google? code quality". In: *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering.* 2022, pp. 1302–1313.

[38]    Anastasia Ruvimova et al. "An exploratory study of productivity perceptions in software teams". In: *Proceedings of the 44th International Conference on Software Engineering.* 2022, pp. 99–111.

[39]    Nils Brede Moe et al. "Improving productivity through corporate hackathons: A multiple case study of two large-scale agile organizations". In: *arXiv preprint arXiv:2112.05528* (2021).

[40]    Margaret-Anne Storey, Brian Houck, and Thomas Zimmermann. "How developers and managers define and trade productivity for quality". In: *Proceedings of the 15th International Conference on Cooperative and Human Aspects of Software Engineering.* 2022, pp. 26–35.

[41]    Paloma Guenes et al. "Impostor phenomenon in software engineers". In: *Proceedings of the 46th International Conference on Software Engineering: Software Engineering in Society.* 2024, pp. 96–106.

[42]    Darja Šmite et al. "From forced Working-From-Home to voluntary working-from-anywhere: Two revolutions in telework". In: *Journal of Systems and Software* 195 (2023), p. 111509.

[43]    Staffs Keele et al. *Guidelines for performing systematic literature reviews in software engineering.* Tech. rep. Technical report, ver. 2.3 ebse technical report. ebse, 2007.

[44]    Andrew MacFarlane, Tony Russell-Rose, and Farhad Shokraneh. "Search strategy formulation for systematic reviews: Issues, challenges and opportunities". In: *Intelligent Systems with Applications* 15 (2022), p. 200091.

[45]    He Zhang, Muhammad Ali Babar, and Paolo Tell. "Identifying relevant studies in software engineering". In: *Information and Software Technology* 53.6 (2011), pp. 625–637.

[46]    Xuetao Li et al. "Systematic literature review of commercial participation in open source software". In: *ACM Transactions on Software Engineering and Methodology* 34.2 (2025), pp. 1–31.

[47]    Joonas Hämäläinen, Teerath Das, and Tommi Mikkonen. "A Systematic Literature Review of Multi-Label Learning in Software Engineering". In: *ACM Transactions on Software Engineering and Methodology* (2024).

[48]    Sin Kit Lo et al. "A systematic literature review on federated machine learning: From a software engineering perspective". In: *ACM Computing Surveys (CSUR)* 54.5 (2021), pp. 1–39.

[49]    Matthew J Page et al. "The PRISMA 2020 statement: an updated guideline for reporting systematic reviews". In: *BMJ* 372 (2021). DOI: 10.1136/bmj.n71. eprint: https://www.bmj.com/content/372/bmj.n71.full.pdf. URL: https://www.bmj.com/content/372/bmj.n71.

[50] Klaas-Jan Stol and Brian Fitzgerald. "The ABC of software engineering research". In: *ACM Transactions on Software Engineering and Methodology (TOSEM)* 27.3 (2018), pp. 1–51.

[51] Robert L. Glass, Iris Vessey, and Venkataraman Ramesh. "Research in software engineering: an analysis of the literature". In: *Information and Software technology* 44.8 (2002), pp. 491–506.

[52] H Rex Hartson, Terence S Andre, and Robert C Williges. "Criteria for evaluating usability evaluation methods". In: *International journal of human-computer interaction* 13.4 (2001), pp. 373–410.

[53] Sandra G Hart and Lowell E Staveland. "Development of NASA-TLX (Task Load Index): Results of empirical and theoretical research". In: *Advances in psychology.* Vol. 52. Elsevier, 1988, pp. 139–183.

[54] Patrícia Silva. "Davis' technology acceptance model (TAM)(1989)". In: *Information seeking behavior and technology adoption: Theories and trends* (2015), pp. 205–219.

[55] Deborah R Compeau and Christopher A Higgins. "Computer self-efficacy: Development of a measure and initial test". In: *MIS quarterly* (1995), pp. 189–211.

[56] Igor Steinmacher et al. "Overcoming open source project entry barriers with a portal for newcomers". In: *Proceedings of the 38th International Conference on Software Engineering.* 2016, pp. 273–284.

[57] James A Russell. "A circumplex model of affect." In: *Journal of personality and social psychology* 39.6 (1980), p. 1161.

[58] Lucian José Gonçales, Kleinner Farias, and Bruno C da Silva. "Measuring the cognitive load of software developers: An extended Systematic Mapping Study". In: *Information and Software Technology* 136 (2021), p. 106563.

[59] Brittany Johnson, Thomas Zimmermann, and Christian Bird. "The effect of work environments on productivity and satisfaction of software engineers". In: *IEEE Transactions on Software Engineering* 47.4 (2019), pp. 736–757.

[60] Prem Devanbu et al. "Analytical and empirical evaluation of software reuse metrics". In: *Proceedings of IEEE 18th International Conference on Software Engineering.* IEEE. 1996, pp. 189–199.

[61] Shraddha Barke, Michael B James, and Nadia Polikarpova. "Grounded copilot: How programmers interact with code-generating models". In: *Proceedings of the ACM on Programming Languages* 7.OOPSLA1 (2023), pp. 85–111.

[62] Raja Parasuraman and Dietrich H Manzey. "Complacency and bias in human use of automation: An attentional integration". In: *Human factors* 52.3 (2010), pp. 381–410.

[63] Samarth Sikand et al. "How much SPACE do metrics have in GenAI assisted software development?" In: *Proceedings of the 17th Innovations in Software Engineering Conference.* 2024, pp. 1–5.

[64] Marshall McLuhan. "Laws of the Media". In: *ETC: A Review of General Semantics* 34.2 (1977). Retrieved from JSTOR, pp. 173–179. URL: http://www.jstor.org/stable/42575246.

[65] Amy J Ko. "Why we should not measure productivity". In: *Rethinking Productivity in Software Engineering* (2019), pp. 21–26.

[66] Daniel Fontanet Losquiño and Tomas Urdell. "Why do developers struggle with documentation while excelling at programming". B.S. thesis. Universitat Politècnica de Catalunya, 2014.

[67] Auste Simkute et al. "Ironies of generative AI: understanding and mitigating productivity loss in Human-AI interaction". In: *International Journal of Human–Computer Interaction* 41.5 (2025), pp. 2898–2919.

## PRIMARY STUDIES

[PS1]    Margaret-Anne Storey and Alexey Zagalsky. "Disrupting developer productivity one bot at a time". In: *Proceedings of the 2016 24th ACM SIGSOFT international symposium on foundations of software engineering*. 2016, pp. 928–931.

[PS2]    Frank F Xu, Bogdan Vasilescu, and Graham Neubig. "In-ide code generation from natural language: Promise and challenges". In: *ACM Transactions on Software Engineering and Methodology (TOSEM)* 31.2 (2022), pp. 1–47.

[PS3]    Justin D Weisz et al. "Better together? an evaluation of ai-supported code translation". In: *Proceedings of the 27th International Conference on Intelligent User Interfaces*. 2022, pp. 369–391.

[PS4]    Sandeep Kaur Kuttal et al. "Trade-offs for substituting a human with an agent in a pair programming context: the good, the bad, and the ugly". In: *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 2021, pp. 1–20.

[PS5]    Ketai Qiu et al. "From Today's Code to Tomorrow's Symphony: The AI Transformation of Developer's Routine by 2030". In: *arXiv preprint arXiv:2405.12731* (2024).

[PS6]    Daniel Russo. "Navigating the complexity of generative ai adoption in software engineering". In: *ACM Transactions on Software Engineering and Methodology* 33.5 (2024), pp. 1–50.

[PS7]    Zhensu Sun et al. "Don't Complete It! Preventing Unhelpful Code Completion for Productive and Sustainable Neural Code Completion Systems". In: *2023 IEEE/ACM 45th International Conference on Software Engineering: Companion Proceedings (ICSE-Companion)*. IEEE. 2023, pp. 324–325.

[PS8]    Silvia Abrahão et al. "Software Engineering by and for Humans in an AI Era". In: *ACM Transactions on Software Engineering and Methodology* (2025).

[PS9]    Jenny T Liang, Chenyang Yang, and Brad A Myers. "A large-scale survey on the usability of ai programming assistants: Successes and challenges". In: *Proceedings of the 46th IEEE/ACM international conference on software engineering*. 2024, pp. 1–13.

[PS10]   Rudrajit Choudhuri et al. "How far are we? the triumphs and trials of generative ai in learning software engineering". In: *Proceedings of the IEEE/ACM 46th international conference on software engineering*. 2024, pp. 1–13.

[PS11]   Daye Nam et al. "Using an llm to help with code understanding". In: *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering*. 2024, pp. 1–13.

[PS12]   Vijayaraghavan Murali et al. "AI-assisted Code Authoring at Scale: Fine-tuning, deploying, and mixed methods evaluation". In: *Proceedings of the ACM on Software Engineering* 1.FSE (2024), pp. 1066–1085.

[PS13]   Shaokang Jiang and Michael Coblenz. "An Analysis of the Costs and Benefits of Autocomplete in IDEs". In: *Proceedings of the ACM on Software Engineering* 1.FSE (2024), pp. 1284–1306.

[PS14]   Ranim Khojah et al. "Beyond code generation: An observational study of chatgpt usage in software engineering practice". In: *Proceedings of the ACM on Software Engineering* 1.FSE (2024), pp. 1819–1840.

[PS15]   Wei Wang et al. "Rocks coding, not development: A human-centric, experimental evaluation of LLM-supported SE tasks". In: *Proceedings of the ACM on Software Engineering* 1.FSE (2024), pp. 699–721.

[PS16]   Albert Ziegler et al. "Productivity assessment of neural code completion". In: *Proceedings of the 6th ACM SIGPLAN International Symposium on Machine Programming*. 2022, pp. 21–29.

[PS17]   Nicole Davila et al. "An industry case study on adoption of ai-based programming assistants". In: *Proceedings of the 46th International Conference on Software Engineering: Software Engineering in Practice*. 2024, pp. 92–102.

[PS18]   Priyam Sahoo et al. "Ansible Lightspeed: A Code Generation Service for IT Automation". In: *Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering*. 2024, pp. 2148–2158.

[PS19]   Jinrun Liu et al. "Chatgpt vs. stack overflow: An exploratory comparison of programming assistance tools". In: *2023 IEEE 23rd International Conference on Software Quality, Reliability, and Security Companion (QRS-C)*. IEEE. 2023, pp. 364–373.

[PS20]   Mohammad Amin Kuhail et al. ""Will I be replaced?" Assessing ChatGPT's effect on software development and programmer perceptions of AI tools". In: *Science of Computer Programming* 235 (2024), p. 103111.

[PS21]   Steven L Tanimoto. "Five Futures with AI Coding Agents". In: *Companion Proceedings of the 7th International Conference on the Art, Science, and Engineering of Programming*. 2023, pp. 32–38.

[PS22]   Simone Mezzaro, Alessio Gambi, and Gordon Fraser. "An empirical study on how large language models impact software testing learning". In: *Proceedings of the 28th International Conference on Evaluation and Assessment in Software Engineering*. 2024, pp. 555–564.

[PS23]   Hussein Mozannar et al. "Reading between the lines: Modeling user behavior and costs in AI-assisted programming". In: *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. 2024, pp. 1–16.

[PS24]   Thomas Weber et al. "Significant productivity gains through programming with large language models". In: *Proceedings of the ACM on Human-Computer Interaction* 8.EICS (2024), pp. 1–29.

[PS25]   Helena Vasconcelos et al. "Generation Probabilities Are Not Enough: Uncertainty Highlighting in AI Code Completions". In: *ACM Transactions on Computer-Human Interaction* (2024).

[PS26]   Crystal Qian and James Wexler. "Take it, leave it, or fix it: measuring productivity and trust in human-AI collaboration". In: *Proceedings of the 29th International Conference on Intelligent User Interfaces*. 2024, pp. 370–384.

[PS27]   Gustavo Pinto et al. "Developer Experiences with a Contextualized AI Coding Assistant: Usability, Expectations, and Outcomes". In: *Proceedings of the IEEE/ACM 3rd International Conference on AI Engineering-Software Engineering for AI*. 2024, pp. 81–91.

[PS28]   Mariana Coutinho et al. "The role of generative ai in software development productivity: A pilot case study". In: *Proceedings of the 1st ACM International Conference on AI-Powered Software*. 2024, pp. 131–138.

[PS29]   Tianyi Chen. "The Impact of AI-Pair Programmers on Code Quality and Developer Satisfaction: Evidence from TiMi studio". In: *Proceedings of the 2024 International Conference on Generative Artificial Intelligence and Information Security*. 2024, pp. 201–205.

[PS30]   Jacques Bughin. "The role of firm AI capabilities in generative AI-pair coding". In: *Journal of Decision Systems* (2024), pp. 1–22.

[PS31]   Danie Smit et al. "The impact of GitHub Copilot on developer productivity from a software engineering body of knowledge perspective". In: (2024).

[PS32]   Asha Rajbhoj et al. "Accelerating software development using generative ai: Chatgpt case study". In: *Proceedings of the 17th innovations in software engineering conference*. 2024, pp. 1–11.

[PS33]   Kathrin Komp-Leukkunen. "How ChatGPT shapes the future labour market situation of software engineers: A Finnish Delphi study". In: *Futures* 160 (2024), p. 103382.

[PS34]   Jacques Bughin. "What drives the corporate payoffs of using generative artificial intelligence?" In: *Structural Change and Economic Dynamics* 71 (2024), pp. 658–668.

[PS35]  Wendy Mendes, Samara Souza, and Cleidson De Souza. ""You're on a bicycle with a little motor": Benefits and Challenges of Using AI Code Assistants". In: *Proceedings of the 2024 IEEE/ACM 17th International Conference on Cooperative and Human Aspects of Software Engineering.* 2024, pp. 144–152.

[PS36]  Nicholas Gardella, Raymond Pettit, and Sara L Riggs. "Performance, Workload, Emotion, and Self-Efficacy of Novice Programmers Using AI Code Generation". In: *Proceedings of the 2024 on Innovation and Technology in Computer Science Education V. 1.* 2024, pp. 290–296.

[PS37]  Eduard Frankford et al. "Ai-tutoring in software engineering education". In: *Proceedings of the 46th International Conference on Software Engineering: Software Engineering Education and Training.* 2024, pp. 309–319.