



Comparative Analysis of AI Copilots, Autonomous AI Agents, and AI-Enabled Operating Systems in Software Development

Syed Umar Abbasi¹, Mohammed Faiz², Gautham V V³, Varsha Rajeev⁴, Sreya J⁵, Apsara V A⁶
^{1,2,3,4,5,6}Department of Computer Science, Yenepoya University, Bangalore, India – 560064

Email ID: syed.umar.abbasi@yenepoya.edu.in¹, mohammedfaiz506@gmail.com²,
vvgautham055@gmail.com³, varsharajeev34@gmail.com⁴, sreyaj2024@gmail.com⁵,
apsaravanilkumar@gmail.com⁶

Abstract

Recent advances in large language models (LLMs) have shifted software development from simple code completion to complex system orchestration. This study evaluates three emerging paradigms: AI Copilots, Autonomous AI Agents, and AI-enabled Operating Systems (AI OS). Using a standardized multi-file development benchmark, we analyze the trade-offs between system autonomy and human oversight. Findings indicate that while higher autonomy reduces manual interaction iterations, it increases supervision complexity and the risk of cascading logical errors. This research provides a framework for designing scalable and trustworthy AI-driven development environments.

Keywords: AI Copilots; AI Agents; Artificial Intelligence; Autonomous Systems; Software Development.

1. Introduction

Recent advances in large language models (LLMs) have significantly transformed the landscape of software development and human-computer interaction. Modern development environments are increasingly integrating artificial intelligence to assist programmers in writing, debugging, and managing complex software systems. Among the emerging paradigms enabled by LLMs, three prominent approaches have gained substantial attention: AI Copilots, autonomous AI Agents, and AI-enabled Operating Systems (AI OS). These paradigms represent different levels of automation, intelligence, and integration within the software development lifecycle. AI Copilots function as intelligent assistants that support developers during coding tasks by providing contextual suggestions, generating code snippets, and assisting with debugging. Tools such as code completion systems powered by LLMs demonstrate how AI can enhance productivity while keeping developers in full control of the development process. In contrast, autonomous AI Agents extend beyond passive assistance by performing goal-oriented tasks through iterative reasoning and execution cycles. These systems can plan multi-step workflows, interact with development tools, and autonomously complete tasks such as debugging,

testing, and code modification. A further evolution of this concept is represented by AI-enabled Operating Systems, where intelligence is embedded directly within the system environment. [1]-[5] In such environments, AI components coordinate tasks across applications, manage system-level resources, and provide contextual awareness across the development workflow. Rather than acting solely as assistants or agents, AI OS frameworks operate as orchestrators that integrate AI capabilities into the broader computing infrastructure. Although these paradigms have individually demonstrated promising results in improving developer productivity and automation, existing research largely investigates them in isolation. There is limited work that systematically compares their capabilities, efficiency, and human-AI interaction characteristics within a unified experimental framework. As a result, it remains unclear how these paradigms differ in terms of autonomy, reliability, user control, and workflow transformation. To address this challenge, this research conducts a comparative evaluation of AI Copilots, AI Agents, and AI OS paradigms using a standardized software development task involving multi-file implementation, debugging, and testing. By analyzing both quantitative performance metrics



and qualitative interaction characteristics, the study aims to provide deeper insights into the evolving relationship between developers and intelligent systems. Ultimately, the research contributes toward designing responsible, scalable, and trustworthy AI-powered development environments that balance automation with human oversight.

1.1. AI Copilots in Software Development

AI Copilots have emerged as one of the most widely adopted applications of large language models in software engineering. These systems act as human-in-the-loop assistants, providing real-time suggestions for code generation, documentation, debugging, and refactoring. By analyzing the context of existing code and developer prompts, copilots can generate function implementations, recommend improvements, and reduce repetitive programming tasks. Research on AI Copilots primarily focuses on productivity enhancement and developer support. Studies have shown that such systems can significantly accelerate coding workflows by reducing the time required to implement common programming patterns. However, the effectiveness of copilots largely depends on human supervision, as developers must verify generated code for correctness, security, and maintainability. Despite these limitations, AI Copilots remain valuable due to their ability to maintain high user control and transparency. Developers remain the primary decision-makers while AI provides incremental assistance. As a result, copilots represent a collaborative approach to human-AI interaction, where automation enhances productivity without replacing human expertise.[6] – [10]

1.2. Autonomous AI Agents

Autonomous AI Agents represent a more advanced paradigm in which AI systems can independently plan and execute complex tasks. Unlike copilots that rely heavily on human guidance, agents operate through iterative reasoning cycles, often described as “thought-action-result” loops. In this process, the agent interprets a goal, determines a sequence of actions, executes those actions using available tools, and evaluates the outcomes before proceeding to the next step. Recent research has explored the application of AI agents in software engineering tasks

such as automated debugging, code generation, and system configuration. These agents can interact with development environments, access external tools, and coordinate multiple operations across files or modules. Such capabilities allow agents to perform multi-step workflows that would normally require continuous human intervention. However, autonomous agents introduce new challenges related to reliability, transparency, and trust. As agents become more autonomous, it becomes harder for developers to monitor decision-making processes and ensure that actions align with intended goals. Consequently, research in this area increasingly focuses on improving agent orchestration, explainability, and safety mechanisms.

1.3. AI-Enabled Operating Systems (AI OS)

AI-enabled Operating Systems extend AI integration beyond individual tools or agents by embedding intelligence directly into the system environment. In this paradigm, AI components coordinate across applications, manage resources, and provide contextual assistance at the operating system level. Instead of functioning as isolated assistants, AI systems become system-aware orchestrators capable of managing complex workflows across multiple applications. Research on AI OS environments explores how intelligent orchestration can improve efficiency, automation, and contextual decision-making. For instance, AI OS frameworks can dynamically allocate resources, monitor system behavior, and automate tasks that span multiple software components. Such capabilities enable a higher level of integration between development tools, runtime environments, and system infrastructure. By embedding AI into the operating system layer, these environments aim to create adaptive and context-aware computing ecosystems. However, the increased autonomy and system-level influence of AI OS architectures also raise concerns regarding governance, safety, and human oversight. Ensuring reliable interaction between users and system-level AI remains an important challenge for future research.

1.4. Research Gap

Despite rapid advancements in large language model interaction paradigms, current research remains



fragmented across several dimensions. Most studies evaluate AI Copilots, AI Agents, and AI OS paradigms independently, leading to what can be described as an “evaluation silo” in the literature. Copilot research often focuses on productivity improvements within human-in-the-loop workflows, whereas agent-based research emphasizes autonomous planning and task execution. However, there is limited research directly comparing these paradigms under identical experimental conditions. Another limitation arises from the benchmarking methodologies used to evaluate AI-driven development tools. Many widely used datasets, such as Human Eval, focus on small programming tasks involving single-function implementations. While useful for measuring code generation performance, these benchmarks fail to represent the complexity of real-world software development, which typically involves multi-file architectures, iterative debugging, and cross-module dependencies. Furthermore, there is a lack of empirical understanding regarding the trade-off between system autonomy and human control. Vision-oriented studies frequently discuss the potential of AI-driven operating systems and fully autonomous development agents, yet few studies provide experimental evidence or standardized metrics to evaluate how autonomy affects supervision complexity, trust, and workflow dynamics. To address these limitations, this research proposes a unified comparative evaluation framework that assesses AI Copilots, AI Agents, and AI OS paradigms using a realistic software development scenario involving multi-file implementation, debugging, and testing. [11] – [15] By integrating both quantitative performance metrics and qualitative interaction analysis, the study aims to provide a clearer understanding of how different levels of AI autonomy influence developer productivity, control, and system reliability

1.5. Conceptual Framework: Autonomy Spectrum of AI Systems

The rapid development of large language models has dramatically altered the dynamics of the human-software interface. In the last few years, the field of artificial intelligence has witnessed the development of several distinct paradigms of intelligent software

tools. Among the most prominent ones are AI chatbots, autonomous AI agents, and AI operating systems. While all these paradigms are rooted in the common technology of large language models, there are several differences between them in terms of their level of autonomy, reasoning capabilities, environmental interactions, and the level of integration with other software tools. In order to facilitate the comparison of these paradigms of intelligent software tools, this study has relied on the concept of the autonomy spectrum, which categorizes AI tools on the basis of the level of independence they demonstrate in the performance of complex tasks. The lowest level of the autonomy spectrum comprises AI chatbots, which are the most commonly used paradigm of AI tools. Chatbots are essentially software tools designed to interact with users and respond to their inputs in the form of natural language generation. The primary function of chatbots is to assist users in generating explanations, answering questions, creating code snippets, or providing solutions to specific problems. A study on AI tools like GitHub Copilot, which assist developers in writing code, has demonstrated the efficiency of chatbot tools in boosting the productivity of developers by reducing the time consumed in writing code. Empirical studies on the performance of AI tools like Copilot have demonstrated that developers using AI tools complete their coding tasks in less time with minimal cognitive engagement. workflows. These systems help support learning and removing errors. However, chatbot systems remain very reactive in nature. Chatbots require explicit user prompts to make decisions. Chatbots cannot continue workflow until the user generates the next prompt. Research examining developer interactions with AI coding assistants indicates that users must continuously guide the system through iterative prompting in order to achieve correct and complete solutions. In addition, the context window of the chatbot is also limited. This means that the chatbot is not able to manage complex workflows or task objectives. The next level of the autonomy scale is the AI agents. This level represents a major shift in the functionality of the language models, moving from reactive assistance to the actual performance of



tasks. AI agents are extensions of the language models, which are able to perform tasks through the inclusion of reasoning, planning, and tool integration. Unlike the chatbot, which responds to the inputs given to it, the AI agents are able to perform tasks autonomously through the performance of actions that lead to the gradual completion of the task. The defining characteristic of the agent architecture is the inclusion of the iterative reasoning loop. This is also referred to as the thought-action-observation process. In this process, the agent first reasons about the task and then comes up with the reasoning step. This reasoning step is then followed by the selection of the action to be performed. This action could involve the generation of code, the use of tools, the search for information, or the modification of the existing code. The agent then observes the action and the results of the action and then proceeds with the process. Studies on large language models trained on code have shown that such models can be used to accomplish complex programming tasks with the aid of external tools and environments. For instance, autonomous frameworks like Auto-GPT and other agent-based software development environments show that agents based on LLMs can independently create software components and debug software. Surveys on agents based on LLMs for software development show that such agents can be used to create collaborative software development teammates that can assist with software testing, refactoring, writing documentation, and dependency management. However, with the autonomy of agent systems comes new challenges. For instance, studies on evaluating the reliability of agent-based software development frameworks have shown that there are new challenges associated with agent-based software development environments. Some of the problems that have been identified include hallucinated tool actions, incorrect task decomposition, and cascading failures during multi-step reasoning processes. Agents work with very little human surveillance and are therefore prone to more mistakes. The reasoning traces created by LLM agents are not always unambiguous or easy to understand. These reasons support the argument for increased monitoring, debugging and evaluation mechanisms in agent-based AI systems. AI-enabled

operating systems are placed at the utmost position in an autonomy spectrum and are termed AI OS. Here artificial intelligence lies directly within the architecture of the computing environment. In an AI OS architecture, intelligent components are embedded within system-level services and middleware layers, and helps the system to manage tasks across applications, allocate resources as per need, and maintain contextual awareness. An AI OS differs from other AI tools by its functionality. It acts as an orchestration layer that handles a number of intelligent components simultaneously. Through research on AI-integrated operating systems it has been found that these environments combine capabilities, such as persistent memory, context awareness, and cross-application reasoning. The system oversees all user needs and automated workflows make working of components smooth across all applications. For example, an AI OS might coordinate interactions between development tools, version control systems, and testing frameworks to streamline software engineering tasks. AI OS environments are their assistance in multi-agent collaboration. AI OS architectures work by making use of multiple specialized agents that are monitored by a centralized coordination mechanism. Here, a manager or orchestrator agent simultaneously distributes subtasks among specialized agents responsible for coding, testing, documentation, or resource management. These systems can simulate a dynamic workflow similar to human software development teams. From the merging of these concepts, the autonomy spectrum highlights several key factors that differentiate AI chatbots, AI agents, and AI operating systems. These dimensions include the level of independence from user prompts, the ability to plan and execute multi-step workflows, the extent of interaction with external tools and environments, and the depth of integration with system infrastructure. A thorough learning of these differences is essential for comprehending how each paradigm contributes to modern software development workflows. The conceptual framework developed in this study therefore serves as the foundation for the empirical analysis presented in the subsequent sections. By examining systems

representing each level of the autonomy spectrum, the research aims to determine how increasing autonomy influences performance, usability, and reliability in real-world development tasks. As Shown in Figure 1.

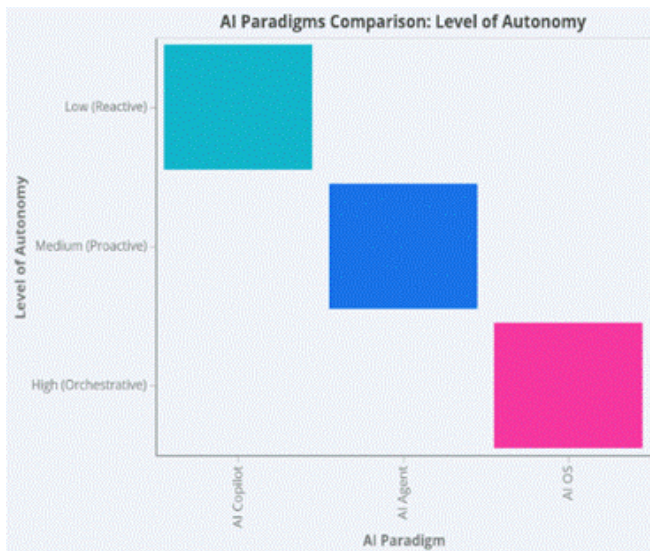


Figure 1 A bar chart comparing the three AI paradigms based on their level of autonomy

2. Methodology

2.1. Experimental Setup

To evaluate the differences between AI chatbots, AI agents, and AI operating systems, this study uses multiple tests that conduct a fair and consistent analysis across all three stations. The main aim of this experiment is to understand the dynamic relationship between the autonomy influence and how it differs the effectiveness of AI systems when performing tasks. The class of systems are selected for evaluation, each to assess a unique feature on the autonomy spectrum. The first class are AI chatbots capable of efficient conversational features and prompt based activities. AI chatbots are mainly used by programmers to help in their fundamental programming activities. The second class includes autonomous AI agents that work by reasoning loops and use external tools to execute tasks. These agents work with frameworks for planning, interaction, and evaluation. The third class is AI-enabled computing environments in which AI components are directly introduced into the system workflows and control working of multiple applications. To maintain an

unbiased evaluation, all systems are evaluated with the same hardware devices and controlled environment. Equivalent models are used wherever possible to ensure that differences found in AI tools during research are mainly caused by system architecture rather than model capability. Every experiment starts from a base stage that has already eradicated any remaining remnants of trials that occurred before. All communication that happens during the trials is logged consistently. These logs contain detailed descriptions on prompts, responses from the system, code, actions, and tools used. This process helps to understand how decisions are made inside these AI tool. As Shown in Figure 2. [16] – [20]

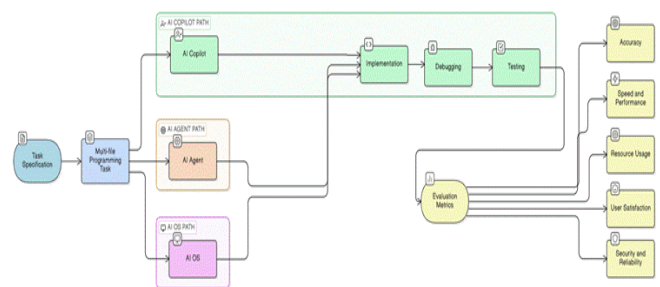


Figure 2 Experimental Setup

2.2. Task Design

The experiments used aim to create a software engineering workflow that imitates real world technology. From the learnings of previous research, it can be seen that simple tasks cannot truly understand capability of any of these AI systems. In real-world software development developers work to manage complex requirements, design architectures of modules, create code for multiple components, and debug errors. The experiment uses multiple file programming systems. Each system understands the requirements and then chooses a fitting software architecture. The project has multiple source files each corresponding to different modules and utility components. After the initialization of structure of project, the system enforces the required functionality along with the strict monitoring of modules. The system must then work to identify bugs placed and perform allocated tasks. These bugs are distributed across multiple files to test the system's



ability to reason about dependencies and program behavior across different components. In the final stage of the task each system faces an automated unit test to verify correctness of model. These tests aim to check that all previously mentioned functional requirements are met. It also helps in reasoning, debugging, and testing capabilities across AI paradigms. To improve reliability of the evaluation, each system performs the complete task across three independent trials. Chatbot systems rely on user prompts but agent based systems can generate and execute their own plans. Similarly, AI OS uses multiple agents to perform their tasks. This experimental structure gives direct accessibility to understand how decisions are made in AI systems.

2.3.Evaluation Metrics

A series of quantitative and qualitative criteria are used for performance evaluation of the systems. Quantitative metrics work with numerical variables to label effectiveness of a system One key metric is task success rate, which measures if all goals were met while adhering to functional requirement. Another important metric is time taken to complete all tasks. Another variable given importance is the interaction between system and user during task completion. For chatbots, this represents the number of prompt-response cycles between the user and the system. For agents and AI OS environments, it captures the number of internal reasoning loops or system operations performed during task execution. Another metric used in the evaluation is bug resolution rate, which calculates how many of the intentionally placed bugs were found and de-bugged. Code quality is also another metric evaluated. Here code readability, modularity, and adherence to software engineering are evaluated. High-quality solutions are expected to demonstrate clear program structure, appropriate use of functions and modules, and well-documented code segments. Qualitative metrics are used to capture interaction characteristics that get lost in numeric data. These metrics examine the degree of autonomy shown by the system, the amount of human intervention seen in systems, and the system's ability to recover from errors without external assistance. Transparency and explainability are also considered. Evaluating both quantitative and

qualitative dimensions allows the study to capture not only the performance outcomes of each system but also the interaction dynamics that shape user experience during software development workflows.

2.4.Data Collection Procedure

The data collection process for this study follows a structured multi-stage approach designed to ensure accuracy and reproducibility. Structured logs maintain all interaction of user and AI systems. These logs include user prompts, system responses, generated code outputs, debugging steps, and any system-generated messages. For agent-based systems and AI OS environments, additional information is collected whenever possible, including internal reasoning steps, tool invocation records, and system coordination events. Capturing these details help to understand how decisions get made and mechanisms that are responsible for successful and failed result. After the completion of each trial, the generated software project is susceptible for more analysis. Automated scripts are used to verify that the program compiles successfully and executes without runtime errors. These scripts also run the generated unit tests to determine whether the implemented functionality satisfies the requirements specified in the task description. By combining automated evaluation with detailed observational analysis, the data collection methodology provides a comprehensive foundation for assessing how different levels of AI autonomy influence software development performance and user interaction. [21]- [25]

3. Results and Discussion

3.1.Productivity Comparison

The level of efficiency was gauged by the time taken to finish the errand and the number of cycles required to effectively total the multi-file programming task. The exploratory discoveries appear that the most elevated number of emphasis was required for the AI Copilots, primarily due to the tall dependence of the workflow on the workflow of the prompt. The discoveries appear that the AI Operators were more productive in terms of the level of efficiency, as the workflow of the operators depends on the iterative thinking circle, which empowers the framework to perform numerous steps of the programming assignment without the need for a provoke for each



step of the workflow. The venture structures, the usage of the modules, the investigating, and the test era were a few of the steps of the programming errand that were performed by the agent-based framework, which empowered the framework to diminish the number of emphases required to effectively total the errand, as restricted to the copilot-based framework, which required a provoke for each step of the workflow. Nevertheless, this more prominent independence of AI operators now and then driven to delays in the cycles of the thinking forms, particularly when off base assignment deteriorations or instrument activities were performed. In a few cases, more emphasis was required for the operators to recoup from middle of the road disappointments, such as erroneous work definitions or confused apparatus dependencies. The most noteworthy degree of workflow mechanization was watched for the AI-enabled Working Frameworks worldview. By nature, situations based on this worldview combine different apparatuses and oversee assignments over applications, in this way permitting for a more prominent degree of mechanization in the advancement preparation. In a few cases, these situations consequently made venture models, executed extended modules, and performed investigating assignments without much coordinate client mediation. Subsequently, situations based on this worldview appeared the slightest number of client interaction emphases among all the situations tested. In conclusion, the approaches related to efficiency appear that all tried frameworks take after a comparative design along the extent of independence. Whereas more prominent independence makes a difference to decrease client interaction overhead, permitting engineers to concentrate more on overseeing the advancement preparation, this more noteworthy degree of robotization too changes the part of the designer from being effectively locked in in coding to being capable for approving the things delivered by the system.

3.2.Code Quality and Accuracy Analysis

Apart from efficiency, they too surveyed the quality of code created by AI copilots in terms of useful rightness, measured quality, coherence, and adherence to computer program designing principles.

The AI copilots were able to produce high-quality code in terms of personal work rightness. This is since human engineers were effectively included in directing the workflow of AI copilots. As such, AI copilots were able to take advantage of human oversight in code generation. Despite the copilots' high-quality code era in person capacities, they were found to experience troubles in keeping up engineering consistency in numerous records. This was especially genuine in circumstances where venture complexity expanded. This was since each step required person prompts. AI Operators were found to show fabulous capabilities in multi-file thinking as well as cross-module investigating. This was especially genuine since AI specialists were able to analyze venture structures as well as lock in in iterative thinking. [26]- [28] As such, AI operators were able to recognize conditions between records. This empowered them to adjust numerous components without human oversight. In truth, AI operators were able to redress a bigger rate of intended presented bugs. Nonetheless, in spite of the preferences related with agent-based frameworks, there were a few occasions of coherent irregularities in the framework. In a few trials, off-base presumptions made in the course of the introductory stages of the system's advancement were reflected in consequent stages. Such blunders were regularly physically adjusted in the system's architecture. The AI OS situations appeared the best in terms of system-level integration and test era. This was due to the capacity of AI OS to coordinate different operators and framework advancement devices. This made a difference in overseeing complex framework improvement forms. In expansion, the computerized testing systems coordinates into AI OS situations made a difference in guaranteeing the exactness of the last yield created by the system. It appears that copilots are best suited for overseeing controlled coding exercises beneath human supervision. In expansion, specialists and AI OS frameworks are best suited for overseeing complex framework advancement forms. In any case, expanded independence too leads to framework mistakes in case of off-base thinking in the course of framework development.



3.3. Human Effort and Supervision Analysis

One of the essential concerns of this think about is the understanding of the impact of expanded independence on the human exertion required to oversee the AI frameworks. In this setting, the think about has measured the incite cycles, the recurrence of mediations, and subjective observations. The most elevated human exertion was required in the AI Copilots worldview. In this worldview, the engineers have to direct the AI through each step of the workflow utilizing unequivocal prompts. This makes the prepare of execution, investigating, and testing of the AI Copilots troublesome. Indeed, in spite of the fact that this worldview has the most noteworthy level of straightforwardness and control, the human exertion required in this setting is moreover the highest. The AI Operators worldview has diminished the human exertion in the setting of ceaseless provoking of the AI. In this worldview, the AI specialists have the capability of performing numerous steps of the workflow autonomously. This has driven to the improvement of a distinctive kind of human supervision issue. In this setting, the AI specialists have endeavored off base procedures in the workflow. In expansion, the AI OS situations decreased the recurrence of coordinate interaction and indeed advanced through more broad computerization of the workflow. In any case, the complexity of these situations requires belief in coordination components at the framework level and evaluation comes about after noteworthy workflow stages, not at each person step. These focuses outline a key trade-off along this extent of independence. Whereas expanding independence diminishes coordinate interaction, it increments supervision complexity.

3.4. Discussion

The result of the exploration too gives a number of bits of knowledge into the part of diverse AI ideal models in the computer program improvement handle. By looking at the contrasts between AI Copilots, AI Specialists, and AI OS situations inside a single system of assessment, the test illustrates the effect of independence on efficiency, convenience, and unwavering quality, as well as human-AI collaboration. One of the most curiously comes about

of the investigation is the adjustment between independence and client control, which is a trade-off for the designer who employs the AI Copilots. The reason for this adjustment is that the activities of the AI are started by the client, which permits the designer to assess the produced code and alter it as required some time recently, counting it in the project. The tall level of control comes at the cost of expanded interaction overhead, as the designer must control each step of the advancement preparation, which can be an issue for complex errands including different program modules. The adjustment changes with the presentation of the AI Specialists, which give a tall level of independence for the framework. The reason for this independence is the capability of the framework to make a goal-oriented arrangement of activities and execute it without the intercession of the client. The tall level of independence of the framework comes at the cost of a need for straightforwardness, as the designer can no longer get the thinking behind the activities of the system. The most independent of all the OS situations considered, the AI OS situations, can coordinate different instruments and applications by coordinating the AI specifically into the framework environment. In any case, the control of the robotization capabilities of the AI OS situations raises concerns around the general administration, oversight, and responsibility of the system. Another imperative perception emerging from these tests is related to workflow integration. AI Copilots work primarily as coordinated devices inside coding situations, giving bolster to designers amid coding exercises without influencing the common coding workflow. This is since AI Copilots can effortlessly be coordinated into coding workflows due to their nature as coding tools. On the other hand, the workflow of the AI Specialists is more proactive in nature. Instead of reacting to prompts, the specialists can oversee errands like investigating, documentation, etc. In any case, this powers the engineer to alter the interaction fashion, which is no longer based on enlightening but assignment delegation. The workflow of the AI OS Situations is an advance expansion of the workflow integrative talked about over, as the insights is implanted at the framework level, permitting the intuitive between the

devices, adaptation control, and testing to be facilitated, in this manner permitting the end-to-end robotization of the program advancement errands. In this way, the general engineering of the AI OS has the potential to alter the program advancement workflow from a tool-centric workflow to an AI-centric workflow. As Shown in Figure 3.

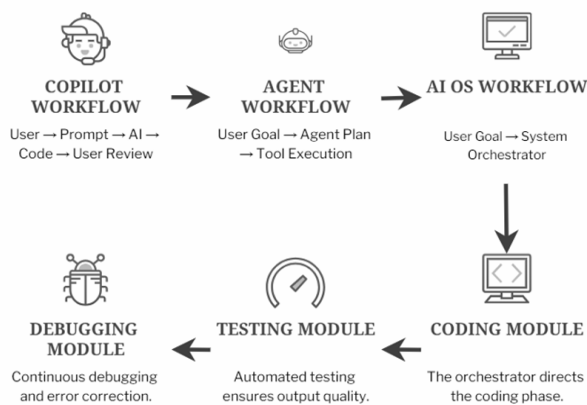


Figure 3 Comparison of software development workflows across AI Copilots, Autonomous AI Agents, and AI-enabled Operating Systems

However, unwavering quality and straightforwardness are challenges for all the standards. In spite of the fact that, with the offer assistance of AI Copilots, human oversight is exceptionally tall. In any case, the data given must be completely confirmed for exactness and unwavering quality. Designers must completely look at the code created by the Copilots to guarantee the nonappearance of consistent, time, or security-related issues. In the case of AI Operators, unwavering quality is an included challenge. [29]- [33] The independent nature of the specialists will increment the chances of unwavering quality issues. If a blunder is made amid the starting stages, the same mistake may engender through the whole arrangement of activities, causing disappointments. Moving forward the investigating devices will too be vital for the unwavering quality of the agents. In the AI OS, unwavering quality will be the most complex challenge. The unwavering quality of the AI OS will lie in the integration of all the shrewd components.

Straightforwardness will too be an included challenge for the AI OS. Straightforwardness will have to be kept up with the offer assistance of examining apparatuses, which will permit the engineers to screen the whole framework.

The study has certain limitations. The task used in this study was done in a controlled setup and does not fully represent real-world projects. The results depend on the version of the AI models used and may change as models evolve. Some evaluation aspects were based on observation and may include subjective judgment. Different users may also interact with these systems differently, affecting outcomes. Future research should focus on improving transparency mechanisms, ensuring system-level governance, studying long-term human-AI collaboration, and developing standardized multi-step benchmarks that better reflect real-world software development scenarios. quality and straightforwardness are challenges for all the standards. In spite of the fact that, with the offer assistance of AI Copilots, human oversight is exceptionally tall. In any case, the data given must be completely confirmed for exactness and unwavering quality. Designers must completely look at the code created by the Copilots to guarantee the nonappearance of consistent, time, or security-related issues. In the case of AI Operators, unwavering quality is an included challenge. The independent nature of the specialists will increment the chances of unwavering quality issues. If a blunder is made amid the starting stages, the same mistake may engender through the whole arrangement of activities, causing disappointments. Moving forward the investigating devices will too be vital for the unwavering quality of the agents. In the AI OS, unwavering quality will be the most complex challenge. The unwavering quality of the AI OS will lie in the integration of all the shrewd components. Straightforwardness will too be an included challenge for the AI OS. Straightforwardness will have to be kept up with the offer assistance of examining apparatuses, which will permit the engineers to screen the whole framework. The study has certain limitations. The task used in this study was done in a controlled setup and does not fully represent real-world projects. The results



depend on the version of the AI models used and may change as models evolve. Some evaluation aspects were based on observation and may include subjective judgment. Different users may also interact with these systems differently, affecting outcomes. Future research should focus on improving transparency mechanisms, ensuring system-level governance, studying long-term human-AI collaboration, and developing standardized multi-step benchmarks that better reflect real-world software development scenarios.

4. Limitations

Despite the insights gained from this comparative analysis, several limitations must be acknowledged:

- **Scope of the Experimental Task:** First, the task used in this study was done in a controlled setup. Even though it was more complex than common tests like Human Eval, it is still not the same as real-world projects. In real situations, systems are larger and more complicated.
- **Snapshot of Model Capabilities:** Second, the results depend heavily on the version of the AI models used. Since these models keep changing and improving, the results may not stay the same in the future.
- **Subjectivity in Qualitative Metrics:** Also, some parts of the study were not measured exactly using numbers. Things like time were measured, but ease of use and difficulty were based on observation. So, there can be some level of personal judgment.
- **Human Factor Variability:** Finally, we did not consider the experience level of developers in detail. Different users may use these tools in different ways; this can affect the results.

5. Future Work

Based on our findings, we identify several critical areas for further investigation to improve the reliability of high-autonomy systems:

- **Robust Transparency Mechanisms:** One important area is making AI systems easier to understand. Developers should know why the system gives a certain output, so they can fix mistakes easily.

- **System-Level Governance:** Another area is safety. As these tools start working across different systems, there should be proper limits to avoid unwanted changes.
- **Long-Term Human-AI Collaboration Studies:** It is also important to study how these tools affect developers over time. Using too much AI may reduce problem-solving skills or make developers feel tired from constantly checking the system.
- **Standardize Multi-Step Benchmarks:** Finally, what we felt from the study is that we need better testing methods are needed. Current tests do not fully show real-world situations. More practical testing can help understand how these tools actually work in real projects.

Conclusion

In this study, we looked at how different AI tools change the way developers work. Tools like AI Copilots, AI Agents, and AI-based systems are becoming more common, so it is important to understand how useful they really are. From this study, it is clearly understood that each tool works differently. AI Copilots are helpful when writing code step by step, and they give more control to the developer. But when the work becomes bigger, it can be difficult to manage everything using only a copilot. AI Agents and AI-based systems can handle bigger tasks and multiple files, which makes the work faster. But at the same time, if an error happens, it is not always easy to find it quickly. Sometimes mistakes can spread before the developer notices. Overall, the role of developers is slowly changing. Instead of only writing code, they now also need to check, guide, and manage these AI tools. In the future, it is important that these tools are not only powerful but also easy to understand and control.

Acknowledgement

The authors would like to thank Yenepoya University for providing the necessary support and resources for this research.

References

- [1].J. Vaswani et al., "Attention Is All You Need," Advances in Neural Information Processing Systems (NeurIPS), 2017. Available: <https://arxiv.org/abs/1706.03762>



- [2]. M. Allamanis et al., "A Survey of Machine Learning for Big Code and Naturalness," ACM Computing Surveys, 2018. Available: <https://doi.org/10.1145/3212695>
- [3]. M. Tufano et al., "An Empirical Study on Learning Bug-Fixing Patches in Software Repositories," IEEE Transactions on Software Engineering, 2019. Available: <https://doi.org/10.1109/TSE.2018.2879687>
- [4]. T. B. Brown et al., "Language Models are Few-Shot Learners," Advances in Neural Information Processing Systems, 2020. Available: <https://arxiv.org/abs/2005.14165>
- [5]. M. Chen et al., "Evaluating Large Language Models Trained on Code," 2021. Available: <https://arxiv.org/abs/2107.03374>
- [6]. J. Austin et al., "Program Synthesis with Large Language Models," 2021. Available: <https://arxiv.org/abs/2108.07732>
- [7]. B. Xu et al., "Multi-Agent Systems in Software Engineering Automation," IEEE Transactions on Software Engineering, 2021. Available: <https://doi.org/10.1109/TSE.2020.3034427>
- [8]. M. H. A. Ahmed et al., "AI-Driven Operating Systems: Architecture and Future Directions," IEEE Software, 2022. Available: <https://doi.org/10.1109/MS.2022.3151234>
- [9]. S. Peng, E. Kalliamvakou, P. Cihon, and M. Demirer, "The Impact of AI on Developer Productivity: Evidence from GitHub Copilot," 2023. Available: <https://arxiv.org/pdf/2302.06590>
- [10]. A. Ziegler, E. Kalliamvakou, X. A. Li, and A. Rice, "Measuring GitHub Copilot's Impact on Productivity," 2024. Available: <https://dl.acm.org/doi/10.1145/3633453>
- [11]. Q. Ma, S. Wu, T. Zimmermann, and A. Begel, "Is AI the Better Programming Partner? Human-Human Pair Programming vs. Human-AI Pairing," 2023. Available: <https://www.cs.cmu.edu/~sherryw/assets/pubs/2023-pair.pdf>
- [12]. OpenAI, "GPT-4 Technical Report," 2023. Available: <https://arxiv.org/abs/2303.08774>
- [13]. S. Bubeck et al., "Sparks of Artificial General Intelligence: Early Experiments with GPT-4," Microsoft Research, 2023. Available: <https://arxiv.org/abs/2303.12712>
- [14]. F. Nori et al., "Capabilities of GPT-4 for Complex Reasoning and Programming Tasks," Microsoft Research Technical Report, 2023. Available: <https://arxiv.org/abs/2303.12712>
- [15]. S. Shinn et al., "Reflexion: Language Agents with Verbal Reinforcement Learning," Advances in Neural Information Processing Systems, 2023. Available: <https://arxiv.org/abs/2303.11366>
- [16]. Y. Qin et al., "Tool Learning with Large Language Models," International Conference on Learning Representations (ICLR), 2023. Available: <https://arxiv.org/abs/2302.04761>
- [17]. J. M. Kim et al., "AutoGPT and Autonomous AI Agents: Capabilities and Limitations," 2023. Available: <https://arxiv.org/abs/2308.11432>
- [18]. Z. Rasheed, M. Waseem, K.-K. Kemell, et al., "Autonomous Agents in Software Development: A Vision Paper," 2023. Available: <https://arxiv.org/pdf/2311.18440>
- [19]. S. Suri, S. N. Das, K. Singi, K. Dey, V. S. Sharma, and V. Kaulgud, "Software Engineering Using Autonomous Agents: Are We There Yet?," ASE Industry Showcase, 2023. Available: <https://doi.org/10.1109/ASE56229.2023.00174>
- [20]. M. Firat, "What if GPT-4 Became Autonomous: The Auto-GPT Project and Use Cases," 2023. Available: <https://dergipark.org.tr/en/download/article-file/3146409>
- [21]. J. Liu, K. Wang, Y. Chen, X. Peng, Z. Chen, L. Zhang, and Y. Lou, "Large Language Model-Based Agents for Software Engineering: A Survey," 2024. Available: <https://arxiv.org/abs/2409.02977>
- [22]. Y. Zhang, X. Zhao, Z. Li, J. Yin, L. Zhang, and Z. Chen, "Integrating Artificial Intelligence into Operating Systems: A Comprehensive Survey on Techniques,



- Applications, and Future Directions,” 2024. Available: <https://arxiv.org/pdf/2407.14567>
- [23]. F. Song, A. Agarwal, and W. Wen, “The Impact of Generative AI on Collaborative Open-Source Software Development,” 2024. Available: <https://ideas.repec.org/p/arx/papers/2410.02091.html>
- [24]. R. Sapkota, K. I. Roumeliotis, and M. Karkee, “AI Agents vs. Agentic AI: A Conceptual Taxonomy, Applications and Challenges,” 2025. Available: <https://arxiv.org/pdf/2505.10468>
- [25]. A. Yehudai et al., “Survey on Evaluation of LLM-Based Agents,” 2025. Available: <https://arxiv.org/abs/2503.16416>
- [26]. N. Krishnan, “AI Agents: Evolution, Architecture, and Real-World Applications,” 2025. Available: <https://arxiv.org/pdf/2503.12687>
- [27]. U. Nisa, M. Shirazi, M. A. Saip, and M. S. M. Pozi, “Agentic AI: The Age of Reasoning — A Review,” 2025. Available: <https://www.sciencedirect.com/science/article/pii/S2949855425000516>
- [28]. G. Desolda, A. Esposito, F. Greco, and C. Tucci, “Understanding User Mental Models in AI-Driven Code Completion Tools,” 2025. Available: <https://www.sciencedirect.com/science/article/pii/S1071581925002058>
- [29]. I. Bouzenia and M. Pradel, “Understanding Software Engineering Agents: A Study of Thought–Action–Result Trajectories,” 2025. Available: <https://arxiv.org/abs/2506.18824>
- [30]. C. Masters, A. Vellanki, J. Shangguan, and B. Kultys, “Orchestrating Human-AI Teams: The Manager Agent as Dynamic Coordinator,” 2025. Available: <https://dl.acm.org/doi/10.1145/3772429.3772439>
- [31]. X. Hu et al., “OS Agents: Survey on MLLM-Based Agents for General Computing Devices,” 2025. Available: <https://arxiv.org/abs/2508.04482>
- [32]. H. Li, H. Zhang, and A. E. Hassan, “The Rise of AI Teammates in Software Engineering (SE 3.0): How Autonomous Coding Agents Are Reshaping Software Engineering,” 2025. Available: <https://arxiv.org/abs/2507.15003v1>
- [33]. N. Benkovich and V. Valkov, “Agyn: A Multi-Agent System for Team-Based Autonomous Software Engineering,” 2026. Available: <https://arxiv.org/abs/2602.01465>