

AI Agentic Programming: A Survey of Techniques, Challenges, and Opportunities

HUANTING WANG, University of Leeds, UK

JINGZHI GONG, University of Leeds, UK

HUAWEI ZHANG, University of Leeds, UK

JIE XU, University of Leeds, UK

ZHENG WANG, University of Leeds, UK

AI agentic programming is an emerging paradigm where large language model (LLM)-based coding agents autonomously plan, execute, and interact with tools such as compilers, debuggers, and version control systems. Unlike conventional code generation, these agents decompose goals, coordinate multi-step processes, and adapt based on feedback, reshaping software development practices. This survey provides a timely review of the field, introducing a taxonomy of agent behaviors and system architectures and examining relevant techniques for planning, context management, tool integration, execution monitoring, and benchmarking datasets. We highlight challenges of this fast-moving field and discuss opportunities for building reliable, transparent, and collaborative coding agents.

CCS Concepts: • **Software and its engineering** → **Software development techniques**.

Additional Key Words and Phrases: Large Language Models, LLMs, AI Agents, AI Agentic Programming

1 Introduction

The software development paradigm is changing rapidly with the rise of large language models (LLMs) [103]. These models enable artificial intelligence (AI) systems that not only translate natural language descriptions into code snippets [119] but also understand task requirements, interact with development tools, and iteratively refine their outputs to produce complex software [42, 191]. Recent studies suggest that developers now use LLMs routinely to assist in daily coding tasks [84, 92, 103]. Unlike traditional code generation tools [41] that respond to a single prompt with a static code snippet, emerging AI coding agents are designed to operate within dynamic software environments, performing iterative, tool-augmented tasks to achieve complex goals.

This shift has given rise to a new programming paradigm, **AI agentic programming**, where LLM-based coding agents can autonomously plan, execute, and refine software development tasks [88, 176]. Unlike conventional code-completion tools [19, 22, 57], which primarily assist with local suggestions, these agents are capable of generating entire programs or modules from natural language specifications, diagnosing and fixing bugs using compiler or test feedback, writing and executing test cases, and refactoring code for readability or performance. Beyond code generation, they can also orchestrate external tools, such as compilers, debuggers, performance profilers, and version control systems, supporting an end-to-end development workflow.

This emerging paradigm has the potential to fundamentally change how software is built and maintained. For example, an AI agent can take a natural language description of a feature and work through a series of steps, such as writing code, generating tests, running those tests, analyzing and fixing issues, and preparing a pull request [191]. Some state-of-the-art coding agents, like Anthropic’s Claude Opus 4 [4], have demonstrated the ability to continue working for hours while maintaining task consistency, avoiding deadlocks, and recovering from failed actions [176, 191].

Authors’ Contact Information: [Huanting Wang](mailto:H.Wang7@leeds.ac.uk), H.Wang7@leeds.ac.uk, University of Leeds, Leeds, UK; [Jingzhi Gong](mailto:J.Gong@leeds.ac.uk), J.Gong@leeds.ac.uk, University of Leeds, Leeds, UK; [Huawei Zhang](mailto:schz@leeds.ac.uk), schz@leeds.ac.uk, University of Leeds, Leeds, UK; [Jie Xu](mailto:J.Xu@leeds.ac.uk), J.Xu@leeds.ac.uk, University of Leeds, Leeds, UK; [Zheng Wang](mailto:z.wang5@leeds.ac.uk), z.wang5@leeds.ac.uk, University of Leeds, Leeds, UK.

These systems can generate and test code, migrate software between frameworks, debug runtime failures, and integrate new features by decomposing complex goals into manageable subtasks [73, 157]. This represents a clear shift from static, one-shot AI-based code generation to *interactive, iterative, and tool-augmented workflows*.

Although progress has been fast, AI agentic programming is still in its early stages. Existing systems vary in architecture, autonomy, tool integration, and reasoning capabilities. There is no standard taxonomy, benchmark suite, or evaluation methodology. At the same time, multiple key challenges remain. These include ensuring reliability and robustness in dynamic environments [103], mitigating errors and hallucinations in generated code [92], extending support beyond dominant languages such as Python to diverse platforms and software ecosystems [237], and embedding safety, trust, and accountability into autonomous behaviours [242].

The success of AI coding agents also depends heavily on their ability to interact effectively with external tools. However, today’s programming languages, compilers, and debuggers are fundamentally human-centric [40, 167]. They are not designed for automated, autonomous systems. These tools often abstract away internal states and decision-making processes to improve usability, ensure portability, and reduce cognitive load for human users [133, 164]. While this abstraction benefits human developers, it may not fit AI agents, which require fine-grained, structured access to internal states, transformation sequences, and validation logic in order to reason about the effects of their actions [55]. Without such access, AI agents struggle to diagnose failures, understand the implications of their changes, or recover from errors in a principled way. For instance, when a code transformation leads to a build failure, the agent needs more than just an error message - it must trace the failure to specific intermediate steps and understand why certain code edits or actions caused the issue. Existing development environments do not provide hooks and feedback mechanisms to support this kind of iterative, tool-integrated reasoning. Similarly, agentic coding systems would benefit from toolchains that support *iterative development*, *state tracking*, and *rich feedback propagation* - capabilities that most conventional tools do not expose. To operate effectively, AI agents may need access to internal compiler representations, transformation traces, symbolic information, and execution metadata.

These challenges show that AI agentic programming is not just a new way of using existing tools. It is a shift that exposes important gaps in how today’s systems software is designed. As the field evolves rapidly, there is an urgent need to clarify its conceptual landscape, identify common patterns and system architectures, and assess the suitability of current development ecosystems. This is the right moment to step back, take stock of recent progress, and lay out the key questions that researchers and developers need to tackle next. Therefore, this survey aims to provide a comprehensive overview of the emerging field of AI agentic programming. Specifically, it covers:

- A conceptual foundation and taxonomy of AI coding agents,
- A review of core system architectures and underlying techniques,
- A summary of the behavior dimensions, operating modes, evaluation strategies and benchmarking practices of AI coding agents,
- A discussion of key challenges and current limitations, and
- An exploration of future research directions, including opportunities to bridge perspectives across disciplines such as programming languages, software engineering, AI, and human-computer interaction.

We focus primarily on **LLM-driven agentic systems for software development**, though many insights extend to general AI agents in other domains like information retrieval [144]. Our goal is to chart the current landscape, clarify foundational concepts, and support the design of robust, efficient, and trustworthy AI agents for programming.

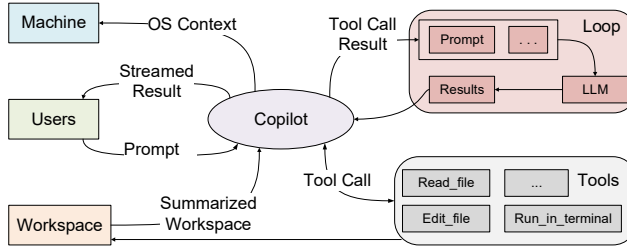


Fig. 1. An example workflow of an AI coding agent.

2 Background

2.1 AI Agentic Programming

AI agentic programming refers to a new programming paradigm in which LLM-based agents autonomously perform software development tasks. Unlike traditional code generation tools that produce outputs in a single step based on a static prompt [237], agentic systems operate in a goal-directed, multi-step manner. They reason about tasks, make decisions, use external tools (such as compilers, debuggers, and test runners), and iteratively refine their outputs based on feedback [150, 176, 205]. These agents can plan sequences of actions, adapt their strategies over time, and coordinate complex development workflows with limited or no human intervention.

At its core, agentic programming combines the capabilities of natural language processing, external tool integration, and task planning. Figure 1 illustrates the architecture of a GitHub Copilot-style agentic programming system [89]. At its core, the agent embeds an LLM within an execution loop, enabling interaction with the development environment. The LLM receives natural language prompts from the user and gathers additional context from the operating system and the workspace (e.g., file summaries or environment state). This information is passed into the reasoning loop, where the LLM decomposes the task into subgoals, generates code or decisions, and determines whether to invoke external tools like reading/editing files or executing terminal commands. Tool outputs are returned to the loop and used as feedback for further refinement. This iterative process continues until the agent completes the task or reaches a stopping condition. Final results are streamed back to the user.

A typical AI agentic programming system is characterized by several key properties. First, it emphasizes *autonomy*, where LLM-based agents can make decisions and take actions without continuous human supervision. Second, it is inherently *interactive*, as agents engage with external tools and environments during execution. Third, it supports *iterative refinement*, allowing agents to improve outputs based on intermediate feedback. Importantly, it is *goal-oriented*, with agents pursuing high-level objectives (e.g., sub-tasks generated from the user inputs) rather than simply responding to one-shot prompts.

Together, these features mark a departure from earlier forms of automation and code generation based on rules [238], classical machine learning models [41] or one-shot LLM calling [237]. AI agentic programming represents a change toward intelligent systems that actively participate in the software development process. This enables new capabilities in intelligent code assistance [89], autonomous debugging and testing [73, 81], automated code maintenance [85], and potentially even self-improving software systems [105, 246].

2.1.1 Working example. As an example of AI agentic programming, consider a developer who tasks an AI coding agent with the following request: “Implement a REST API endpoint that returns the top 10 most frequently accessed URLs from a web server log file. Include unit tests and documentation.”

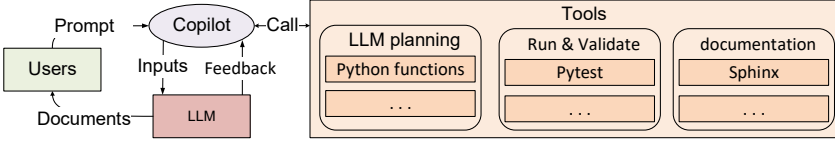


Fig. 2. Agentic workflow for implementing a REST task.

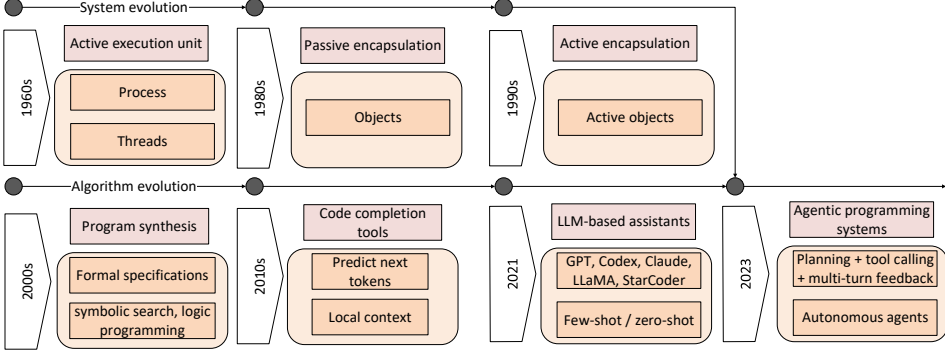


Fig. 3. The evolution of coding agents from program synthesis, to code completion tools, to AI coding agents.

This task requires integrating multiple software components, including file parsing, frequency analysis, web API implementation, testing, and documentation. A high-level view of an agentic loop for solving this task is depicted in Figure 2. Here, an LLM begins by analyzing the natural language task and planning a sequence of actions. It first produces a Python function to parse the log file and count URL frequencies using a dictionary or a data analysis library like `collections.Counter`. Next, it implements a REST API endpoint using a web framework such as `Flask`, exposing a route like `/top-urls` that returns the computed result in JSON format. The LLM then writes unit tests and calls a Python interpreter to execute the generated Python script in the terminal and collects output to validate both the parsing logic and the API usage. It runs these tests using a tool like `pytest`, identifies failing cases, and refines the implementation. If a test fails due to a corner case (e.g., missing fields or malformed input), the LLM goes back to change the Python code, e.g., by adding input validation. It repeats this process - running tools, interpreting results, and modifying the code - until all tests pass. Finally, the LLM can generate documentation strings for each function and call an external tool like `pdoc` or `Sphinx` to produce human-readable API documentation. The process concludes when the agent validates that the tests pass, the API behaves as expected, and the documentation is complete. This example illustrates the core features of agentic programming: autonomous planning, tool integration, iterative refinement, and goal-directed behavior. Unlike one-shot code generation, the agent interacts continuously with tools, learns from feedback, and adapts its actions to deliver a complete and functional software component.

2.2 Historical Context and Motivations

2.2.1 A systems perspective. As illustrated in Figure 3, the idea of automating software development has long been a goal in artificial intelligence and software engineering research [207]. In some respects, the rise of AI agents echoes the historical evolution of operating systems. Early operating systems introduced processes as the basic active execution units, capable of running independently.

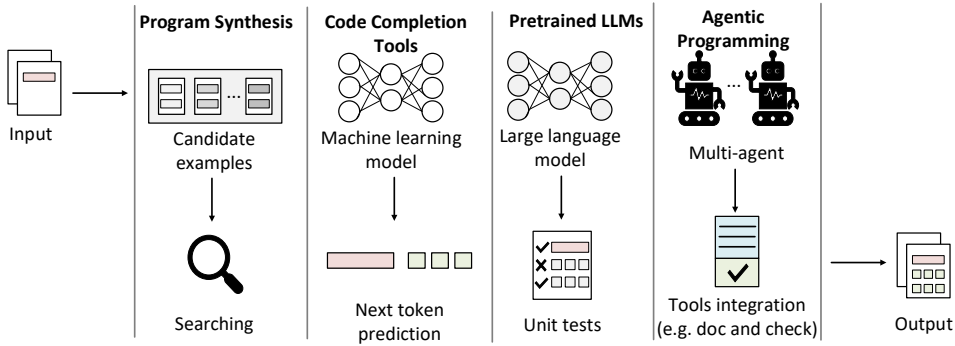


Fig. 4. Evolution of code generation approaches.

These evolved into threads, lightweight abstractions that enabled efficient concurrency and finer-grained scheduling [233]. Later, with object-oriented design [197], objects became the dominant unit of modularity - encapsulating state and behaviour, but remaining fundamentally passive. Active objects [173, 174] extended the object model by incorporating their own thread of control and asynchronous message passing, allowing them to act autonomously rather than merely responding to external calls—an early precursor to modern agents.

LLM-based AI agents can be seen as the next step in this trajectory. Like processes and threads, they are active entities, but unlike objects, they do not wait passively for instructions at every step. Instead, they can plan, act, and coordinate tasks proactively across diverse contexts. Their behavior is not restricted to fixed, handwritten policies but is driven by LLMs, enabling them to adapt dynamically, invoke external tools, and collaborate with both humans and other agents [145]. In this sense, agentic programming represents a new stage in the long-standing pursuit of making software development more autonomous, where code generation, debugging, and optimization are no longer fully hand-crafted but emerge through iterative interactions between humans, tools, and intelligent agents.

2.2.2 Algorithm evolution of agents. Figure 4 shows the evolution of algorithms. Early efforts in *program synthesis* aimed to generate correct-by-construction programs from formal specifications [94], while *code completion tools* sought to improve developer productivity by predicting likely code snippets based on contexts [44, 89, 215]. These approaches, while impactful, typically relied on classical machine-learning models [41], handcrafted rules [238], or statistical techniques with limited generalizations [58].

The advent of large-scale pre-trained LLMs such as Codex [179], StarCoder [160], LLaMas [168], GPTs [180], Claude [49], Gemini [82], Qwen [111], and Deepseek [95], marked a major turning point. These models demonstrated strong zero-shot and few-shot capabilities in generating code [136, 188], translating between programming languages [64, 96], and answering complex programming questions with little or no task-specific fine-tuning [67, 146]. Their ability to understand and generate natural language and code made them a good fit for software development tasks beyond basic code completion [113], including documentation generation [80], test synthesis [123], and bug detection and fixing [73, 81, 125]. As these models became more capable, a new opportunity emerged: using LLMs not just as passive code generators, but as autonomous agents that could reason about goals, invoke tools, and refine their outputs over multiple steps. This shift leads to the paradigm of AI agentic programming, where models operate as task-driven entities capable of planning, interacting with compilers and debuggers, and self-correcting based on feedback.

Several trends motivated this change. First, real-world software development often requires iterative problem solving, tool use, and adaptation, which single-step code generation cannot handle effectively [145, 150]. Second, the rise of prompt engineering and structured prompting techniques (e.g., ReAct, chain-of-thought, scratchpads) enabled LLMs to reason more effectively over multiple steps [95, 243]. Third, the increasing availability of APIs, command-line tools, and language server protocols made it possible to integrate LLMs into full-stack development environments [50, 52, 66]. These developments prompted a rethinking of how LLMs could be deployed—not just as smart coding assistants [89, 215], but as semi-autonomous agents capable of carrying out software engineering tasks with minimal human supervision.

2.3 Agency in AI Systems

Agency is a foundational concept in the design of intelligent systems. At its core, an agent is an entity capable of perceiving its environment, reasoning about goals, and taking actions to influence outcomes. In the context of AI coding agents, *agency* refers to a system’s capacity to act autonomously, i.e., selecting actions based on internal objectives, external feedback, and learned knowledge.

Classical AI research has explored agency extensively in domains such as planning, robotics, and multi-agent systems [145, 235]. Traditional agent models are typically characterized by four key attributes: *reactivity*, the ability to respond to changes in the environment; *proactivity*, the pursuit of long-term goals; *social ability*, the capacity to communicate and coordinate with other agents or humans; and *autonomy*, the ability to operate without direct human intervention. In the context of AI agentic programming, these notions of agency are realized in new ways. An LLM-based system can interpret open-ended tasks expressed in natural language, plan sequences of development steps such as writing, testing, and debugging, and invoke or coordinate external tools like compilers, test runners, and linters. It can also adapt its actions in response to environmental feedback, such as compiler errors or test failures, while maintaining coherent state and reasoning across multiple iterations. Unlike symbolic AI agents that rely on explicitly defined world models and search-based planning, LLM-based coding agents operate in a probabilistic, language-driven manner. Despite this difference, they increasingly exhibit behaviors aligned with classical definitions of agency, especially when augmented with memory, tool-use modules, and planning routines.

2.4 Key Enablers of Agentic Programming

Figure 1 shows a representative system architecture of AI agentic programming. The emergence of AI agentic programming has been made possible by a combination of advances in language modeling, interaction frameworks, and software toolchains. Together, these enablers allow LLMs to move beyond static code generation toward goal-driven, interactive behavior. Below, we summarize the core technical factors that underpin this transition.

2.4.1 Large language models. LLMs trained on massive corpora of code and natural language form the foundation of modern agentic programming systems. These models, as represented by GPT-5 [181], Claude [49], DeepSeek [95], and Gemini [46], serve as the core reasoning engines, powering code generation, task planning, debugging, documentation, and natural language interaction. Their ability to understand and execute complex instructions makes them central to the design of agentic workflows. Modern LLMs can generate syntactically correct and semantically meaningful code, answer development-related queries, and engage in multi-turn conversations with minimal task-specific fine-tuning. Many of these models leverage few-shot, zero-shot, and in-context learning capabilities, allowing them to generalize across programming languages, frameworks, and task domains. This flexibility enables developers to use the same underlying model for a wide range of

Table 1. Representative LLMs for coding tasks.

Model	Size (B)	Context Win.	Tool use	Provider (access)	Open Weight	MoE	Used in coding IDEs
GPT-5	N/A	1 M	✓	OpenAI (API only)	✗	N/A	VS Code, Cursor, other IDEs
GPT-4 variants (o3, o4, etc.)	N/A	128k	✓	OpenAI (API only)	✗	✗	VS Code, JetBrains, Cursor
Claude 4 Opus	~300	200k	✓	Anthropic (API only)	✗	✗	Cursor, Replit (chat)
Gemini 2.5 Pro	~200	1M	✓	Google (API only)	✗	✓	Replit, Google Colab
Grok 4	~1.7T	~128k	✓	xAI	✗	N/A	Not publicly integrated
DeepSeek R1-0528	671 (act. 37)	160k	✓	DeepSeek (API + weights)	✓	✓	Emacs, VS Code (via extension)
Kimi K2	1000 (act. 32)	128k	Limited	Moonshot AI (API)	✓	N/A	Custom plugin support
Qwen3-235B-A22B	235 (act. 22)	128k	Limited	Alibaba	✓	✓	Alibaba Cloud IDE
Qwen3-Coder-480B-A35B-Instruct	480 (act. 35)	256k		Alibaba	✓	✓	Alibaba Cloud IDE
Solar-Pro	72	128k	✓	Upstage (weights on HF)	✓	✗	VS Code (via third-party)
Openhands-LM-32B-v0.1	32	128k	✓	OpenHands	✓	✗	VS Code (via extension)
Devstral-Medium	N/A	128k	✓	Mistral (API only)	✗	N/A	VS Code (via API)
Devstral-Small	24	128k	✓	Mistral (API only)	✓	✗	VS Code (via API)

software engineering tasks, from scaffolding and unit test generation to bug repair and performance tuning. In addition to general-purpose models, some LLMs like Grok [236] and Calude Opus [51] are increasingly optimized for coding tasks through specialized instruction tuning, extended context length, tool use capabilities, and integration with retrieval-based systems. These enhancements make them suitable for multi-turn reasoning, code synthesis grounded in external context, and tool-augmented workflows. Table 1 provides a comparative overview of some of the state-of-the-art LLMs used in code-related tasks. The table compares their key attributes. As the capabilities of LLMs continue to evolve, selecting and fine-tuning the appropriate foundation model for coding tasks becomes a critical design choice in building reliable, efficient, and adaptive agentic systems.

2.4.2 Prompt engineering and reasoning strategies. Effective agentic behavior often requires structured prompting techniques to guide LLMs through multi-step reasoning and tool use. Rather than relying on a single input-output exchange, these methods provide scaffolding that helps models break tasks into manageable steps and maintain coherence across longer interactions. For instance, chain of thought prompting [231] encourages explicit reasoning traces, making intermediate steps visible and improving problem-solving accuracy. ReAct (reasoning and acting) interleaves reasoning with concrete actions [243], such as tool calls or environment interactions, enabling agents to both deliberate and act in context. Scratchpad prompting [38] provides a working memory space where partial results, hypotheses, or plans can be written down and refined, supporting iterative refinement. Modular prompting [214], on the other hand, separates tasks into distinct functional roles—such as planner, executor, and verifier—so that the model can coordinate across specialized subtasks. Together, these techniques allow agents to decompose complex problems, retain intermediate states, and revise their behavior in light of new evidence. They also increase transparency by exposing the reasoning process and provide greater controllability by constraining how models structure their outputs. In practice, structured prompting forms the backbone of many agentic systems, enabling LLMs to move beyond ad hoc responses and toward more reliable, interpretable, and goal-directed behavior.

Table 2. Examples of tools supported by GitHub Copilot agent.

Tool Type	Examples
Compiler	gcc [16], clang [7], javac [20], tsc [34]
Debugger	gdb [14], lldb [23], pdb [28]
Test Framework	pytest [32], unittest [35], Jest [21], Mocha [24]
Linters	eslint [12], flake8 [13], black [3], prettier [30]
Version Control	git [15]
Build System	make [17], cmake [8], npm [26], maven [2]
Package Manager	pip [29] yarn [37], cargo [5]
Language Server	pyright [31], tsserver [33]

```
import OpenAI from "openai";
const client = new OpenAI();
const tools = [
{
  "type": "function",
  "function": {
    "name": "compile_code",
    "description": "...",
    "parameters": {
      "type": "object",
      "properties": {
        "language": {
          "type": "string",
          "enum": ["c", "cpp"],
          "description": "Programming language to compile."
        },
        "source": {
          "type": "string",
          "description": "Source code to compile."
        },
        "flags": {
          "type": "array",
          "items": { "type": "string" },
          "description": "Compiler flags."
        }
      }
    },
    "required": ["language", "source"],
    "additionalProperties": false
  }
},
],
];
```

Listing 1. An example of the OpenAI tool schema

Table 3. Example interfaces between LLMs and tools.

Interface Type	Description	Example
Natural Language	LLM sends plain text instructions; the tool parses heuristically.	"Run my tests and show errors"
Command Line	Tools invoked with shell commands; I/O as text streams.	gcc main.c -O2
Language Server Protocol	JSON-RPC protocol exposing AST, symbols, diagnostics.	VS Code using LSP for Python
REST / gRPC APIs	Tools exposed as network services with structured request/response.	GitHub Actions REST API
Structured Schema (JSON)	Actions, parameters, outputs defined in machine-readable schema.	OpenAI function calling JSON schema
Intermediate Representation	LLM interacts with compiler/runtime IR or AST.	LLVM IR for code optimization
Event Stream / Logs	Tools output execution traces or state changes as structured streams.	Debug logs from pytest
Framework-based Adapters	Middleware unifies heterogeneous tool interfaces.	LangChain, Model Context Protocol

2.4.3 *Tool use and API integration.* Agentic systems rely heavily on external tools, such as compilers, debuggers, test frameworks, linters, and version control systems, to validate and refine generated code. These tools provide the concrete signals needed to check correctness, enforce coding standards,

and ensure that outputs remain consistent with project requirements. For example, Table 2 lists a subset of the tools currently supported by the GitHub Copilot Agent [89], covering compilation, testing, and version control. Integration with external tools can take multiple forms, including command line interfaces, language server protocols (LSP), and RESTful APIs. Increasingly, LLMs interact with tools through structured Python or JavaScript interfaces that specify the available actions, input parameters, and expected outputs in a machine-readable format. For example, Listing 1 shows how a compiler can be exposed as a tool to extend an LLM’s capabilities. This structured approach reduces ambiguity, grounds commands in the correct syntax, and makes it easier for the model to call external tools safely and consistently. By interpreting the schema, an LLM can generate well-formed commands, parse structured responses, and adjust its behavior in a predictable way. Table 3 summarizes the main types of interfaces through which LLMs interact with external tools. These range from free-form natural language instructions to highly structured schemas and domain-specific protocols.

2.4.4 State and context management. LLMs operate under fixed context windows, limiting their ability to reason over long histories. Agentic systems therefore incorporate external memory mechanisms to store plans, results, tool outputs, and partial progress. This memory can take the form of vector stores, scratchpads, or structured logs, allowing the agent to recall relevant information across multiple steps and maintain coherence over long-running tasks. Table 4 compares the context management strategies of mainstream AI coding agents, revealing substantial differences in context size and memory persistence. Tools like GitHub Copilot [89] currently do not utilize persistent memory, instead using transient methods such as sliding windows or dynamic token budgeting. In contrast, agents like SWE-agent [242], Devika [11], and OpenDevin [27] employ persistent storage, often via vector databases or structured stores, to support long-term recall of plans, tool outputs, and project history. Some, such as Cursor IDE [70] and Continue.dev [9], use embedding-based search to retrieve semantically relevant content, while others summarize prior actions to stay within the available context window. These differences reflect a clear trade-off: smaller context windows typically rely on lightweight retrieval or summarization, whereas larger windows with persistent memory enable richer state tracking but add storage and retrieval overhead.

2.4.5 Feedback loops and self-improvement. Agentic programming leverages feedback to refine outputs iteratively. Agents may rerun failed tests, revise prompts based on compiler errors, or reflect on past failures to improve future behavior. For instance, compiler errors may trigger targeted code edits, test failures may prompt iterative debugging, and linter warnings may guide stylistic refinements. Some systems incorporate explicit planning, retry mechanisms, or even gradient-based updates [109] through fine-tuning or reinforcement learning. This closed-loop design supports robustness and adaptability in complex programming tasks.

2.5 Comparison to Related Paradigms

AI agentic programming represents a distinct paradigm that builds upon but fundamentally differs from existing paradigms that have shaped the landscape of automated software development.

2.5.1 Program synthesis. Program synthesis has been a foundational approach to automated code generation, traditionally divided into two types: *deductive synthesis* uses formal specifications to generate provably correct programs, while *inductive synthesis* learns from input-output examples with symbolic search and logic programming techniques to infer program logic [94, 116]. Classical synthesis systems like sketching [207] and more recent neural approaches like RobustFill [76] specialize in generating targeted code snippets that satisfy precise specifications. However, classical program synthesis focuses on single-function generation from formal specifications (due to the

Table 4. Context management mechanisms supported by mainstream AI coding agents.

Agent	Underlying Model	Context Window (default)	Persistent Memory	Context Management Mechanism
GitHub Copilot	GPT-4 (o-series)	16k	✗	Sliding window over active buffer
Codeium	GPT-4 / Claude 3.5	32k	✗	Dynamic token budgeting based on file proximity and edit history
Cursor IDE	Claude 3.5 Sonnet / GPT-4	128k	✓	Semantic search over project history
SWE-agent	GPT-4	16k	✓	Vector DB retrieval for tool outputs and plan state
Devika	GPT-4 / Open LLM	32k	✓	Structured memory via SQLite and embeddings
AutoDev	GPT-4	16k	✓	Summarization of prior actions and tool logs
Continue.dev	GPT-4 / Claude	32k	✗	Embedding-based local recall over recent edits
OpenDevin	GPT-4 / Claude / Mixtral	32k	✓	RAG over command history, plans, and intermediate outputs

scaling challenge of code synthesis [94]) and typically operates in a one-shot generation mode [76], whereas agentic programming handles multi-step workflows (e.g., planning, tool use, and iterative refinement) and engages in continuous interaction with development environments [56].

2.5.2 Code completion tools. Code completion, as one of the most commercially successful applications of AI in programming, excels at context-aware code suggestion, leveraging large-scale pre-training on code repositories to predict next tokens of partially written code [113, 146, 175]. Advanced completion tools can suggest entire functions, classes, or small modules based on comments, function signatures, and surrounding context, with tools like GitHub Copilot [89], TabNine [215], and Amazon Q Developer [44] achieving widespread adoption. These code completion tools often operate as reactive assistants that respond to developer input, while agentic programming systems demonstrate proactive behavior and autonomous planning. Furthermore, agentic programming extends beyond code generation to encompass testing, debugging, deployment, and maintenance activities that completion tools typically do not address [56, 242].

2.5.3 DevOps automation. DevOps automation focuses on streamlining software delivery pipelines through Infrastructure as Code (IaC), Continuous Integration/Continuous Deployment (CI/CD), and automated testing frameworks [112]. Tools like Jenkins [190], GitLab CI [114], and modern platforms like GitHub Actions [90] automate repetitive deployment tasks, testing workflows, and infrastructure management. While both paradigms emphasize automation, DevOps automation primarily handles pre-defined workflows and infrastructure management, whereas agentic programming focuses on adaptive problem-solving and creative solution generation. Additionally, agentic programming can potentially orchestrate and improve DevOps processes themselves, representing a higher-order form of automation [124].

2.5.4 Automated machine learning and automated development. Automated Machine Learning (AutoML) represents a successful paradigm for democratizing AI model development through automation of model selection, hyperparameter tuning, and feature engineering [99]. Platforms like Google Cloud AutoML [65], Amazon SageMaker Autopilot [202], and open-source frameworks like Auto-sklearn automate the traditional machine learning pipeline from data preprocessing to model deployment [86]. However, AutoML focuses on statistical model optimization within well-defined machine learning workflows and operates with structured data and standardized evaluation metrics, while agentic programming tackles more general software development challenges with creative problem-solving and multi-modal reasoning.

2.5.5 Multi-Agent systems and human-AI collaboration. Traditional multi-agent systems in software development typically involve specialized agent roles working within predefined coordination protocols [117]. These systems often feature separate agents for requirements analysis, code generation, testing, and documentation, coordinating through structured communication interfaces. Recent advances in LLM-based multi-agent programming have demonstrated the potential for more sophisticated collaboration, with systems like MetaGPT showing how multiple AI agents can

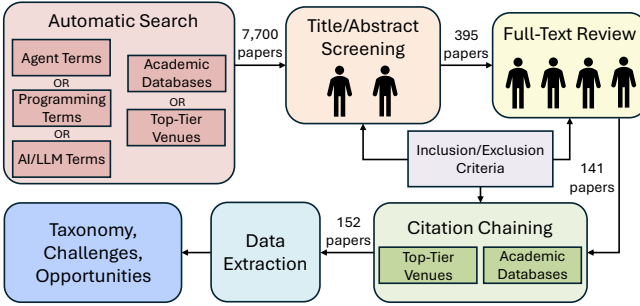


Fig. 5. Survey methodology for academic research.

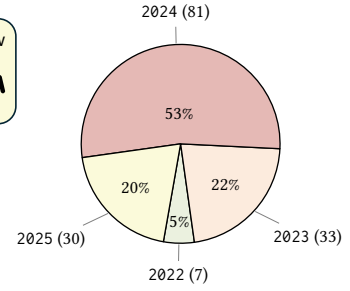


Fig. 6. Distribution of academic papers.

simulate software development teams [100]. AI agentic programming can be viewed as an evolution of multi-agent systems that incorporates human-in-the-loop collaboration and dynamic role adaptation. Unlike traditional multi-agent systems with fixed agent roles and rigid communication protocols, agentic programming systems demonstrate fluid role assignment and context-adaptive behavior. Moreover, the integration of tool use, environmental interaction, and persistent memory distinguishes modern agentic programming from earlier multi-agent approaches. Contemporary agentic systems like AutoGen [235] and CrewAI [79] enable agents to directly interact with development tools, maintain context across extended sessions, and learn from past interactions [145]. This represents a significant advancement over traditional multi-agent systems that typically operated in more constrained, simulation-based environments.

2.5.6 Comparison with robotics and reinforcement learning agents. Robotic agents typically interact with the physical world through sensors and actuators. Their perception, control, and planning modules are tightly coupled with real-time feedback and often require safety guarantees. Reinforcement learning (RL) agents [210], by contrast, learn behaviors by maximizing cumulative reward through trial and error, often in simulated environments. These agents explore large state-action spaces and acquire policies over time. Agentic programming systems share features with both paradigms. Like robotic agents, coding agents must coordinate perception, such as task understanding, with actions such as code edits or tool invocations, all within an environment shaped by constraints and feedback. Like RL agents, they benefit from feedback loops, for example, test results or compiler outputs, and may employ exploration, retry strategies, or even reward-guided behavior. At the same time, coding agents differ in multiple ways. They operate in a symbolic, tool-rich environment where actions are language-based and environments, such as codebases, APIs, and test harnesses, are highly structured. Success requires reasoning not only about immediate feedback but also about abstract software goals, dependencies, and long-term coherence across multiple steps. This makes agency in programming uniquely challenging and distinct from both the physical world and simulated agents.

3 Survey Methodology

This survey follows a widely-used systematic literature review (SLR) methodology [91, 103, 128, 195, 227, 249] to provide comprehensive coverage of AI agentic programming research, as illustrated in Figure 5.

3.1 Search Strategy

We conducted automatic searches across multiple academic databases, including Google Scholar, ACM Digital Library, IEEE Xplore, SpringerLink, and arXiv.org. We also examined proceedings from top-tier venues (FSE, ICSE, ASE, ICML, NeurIPS, AAAI, etc.). As this fast-evolving field is largely dominated by industry, we also pay attention to open-source and the industry releases of relevant results and tools.

Our search string combined the following term clusters using Boolean operators:

- **Agent terms:** “AI agent” OR “agentic” OR “autonomous agent” OR “coding agent” OR “software agent” OR “intelligent agent” OR “task agent” OR “LLM agent”
- **Programming terms:** “programming” OR “coding” OR “software development” OR “code generation” OR “software engineering” OR “developer” OR “autonomous coding” OR “software automation”
- **AI/LLM terms:** “large language model” OR “LLM” OR “language model” OR “foundation model” OR “AI model” OR “neural code generation”

3.2 Study Selection

After initial retrieval, we followed a three-stage study selection process for academic papers: (1) title and abstract screening by two independent researchers, (2) full-text review with disagreement resolution through discussion, and (3) backward and forward citation chaining to identify additional relevant studies. During the selection process, we used the following criteria:

Inclusion criteria - studies were included if they met all of the following:

- (1) Focus on AI systems for software development with autonomous/semi-autonomous behavior
- (2) Demonstrate agentic behaviors: planning, tool use, iterative refinement, or adaptive decision-making
- (3) Present novel techniques, architectures, evaluations, or comprehensive analysis
- (4) Include experimental evaluation, case studies, or substantial implementation details
- (5) Written in English with accessible full text

Exclusion criteria - studies were excluded if they:

- (1) Focused solely on traditional code completion without agentic behavior
- (2) Addressed non-programming domains (robotics, game playing, etc.)

3.3 Results

Our systematic search yielded:

- **Initial retrieval:** 7,700 papers from database searches
- **Title/abstract screening:** 395 papers selected for full-text review
- **Full-text review:** 141 met all criteria.
- **Final corpus:** 152 papers included after full-text evaluation and citation chaining
- **Software tools and industry products:** we also study a wide range of state-of-the-art AI coding agents and LLMs like GitHub Copilot Agents, GPT, Gemini, Deepseek, Qwen and Claude Opus 4.

Among the 152 academic references published from 2022 to 2025 (excluding tool descriptions and websites), 5% appeared in 2022, 22% in 2023, 53% in 2024, and 20% in 2025, as shown in Figure 6, reflecting a surge in AI agent programming research following the widespread adoption of LLMs. Based on our systematic analysis, we developed a hierarchical classification of AI agentic

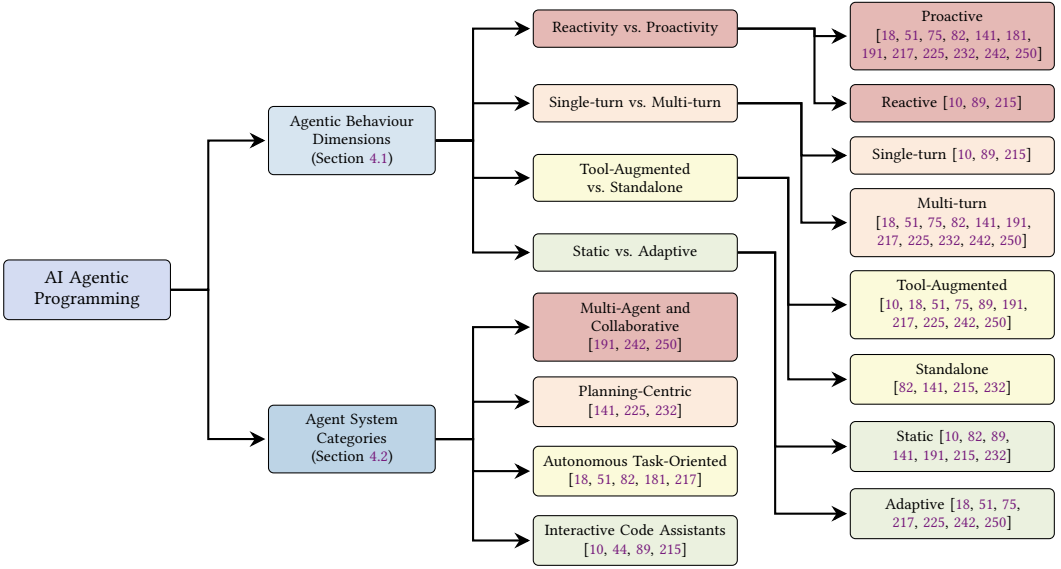


Fig. 7. Taxonomy of AI agentic programming systems.

programming systems along behavioral characteristics and system architectures, as shown in Figure 7. The following section examines each category in detail.

4 Taxonomy of AI Agentic Programming

AI agentic programming is an emerging paradigm that equips LLM-based systems with autonomy, enabling them to plan, execute, and refine programming tasks over multiple steps. To provide structure to the diverse and fast-evolving landscape of agentic programming systems, this section introduces a taxonomy based on key behavioral and architectural dimensions. We categorize existing systems and approaches along these axes to clarify the design space and inform future development.

4.1 Agentic Behaviour Dimensions

This subsection defines the primary behavioural traits that differentiate agentic systems, forming the basis for a comparative classification.

4.1.1 Reactivity vs. proactivity. Reactive agents respond directly to user prompts or feedback without independent task planning. For example, GitHub Copilot reacts by instantly suggesting a function body based on context after users type a function header like `def initial`. It does not include subsequent steps like writing tests or checking generated code. Proactive agents initiate sub-tasks, form execution plans, and re-evaluate decisions, often working autonomously over extended periods. For example, a proactive LLM-based agent can decompose a task into subtasks and execute them automatically. Users providing high-level instructions, such as “*Add a user authentication module*,” will generate the login logic, update the database, integrate it with the UI, and write corresponding unit tests.

4.1.2 Single-turn vs. multi-turn execution. Single-turn agents perform actions in response to individual prompts, often without preserving context. For example, classical GitHub Copilot responds to each prompt independently, without remembering past interactions. In contrast, multi-turn agents,

such as GitHub Copilot Agent or Claude Opus 4 with tooling capabilities, maintain state across interactions, enabling iterative refinement, exploration, and goal pursuit. For instance, GitHub Copilot Agent can hold a conversation across multiple steps, remember earlier function names, and build a complete module through back-and-forth iterations between agents [89].

4.1.3 Tool-augmented vs. standalone agents. Some agents are tightly integrated with external tools (e.g., compilers, debuggers, browsers, test frameworks), allowing them to perform code execution, validation, and correction. Others operate solely within the LLM’s reasoning capabilities, limiting their interactivity and adaptability.

4.1.4 Static vs. adaptive agents. Static agents (e.g., GitHub Copilot and Tabnine [215]) follow predefined workflows or heuristics. Adaptive agents modify their strategies using feedback from tools, user input, or environmental signals. Some employ learning mechanisms to improve over time. For example, GitHub Copilot Agent adapts its approach when test failures occur, revising its implementation or replanning subtasks.

4.2 Agent System Categories

We now present a classification of current AI agentic programming systems, organised by their core functionality and architectural patterns.

4.2.1 Interactive code assistants. These are among the most widely adopted applications of LLMs in software development. These systems assist developers by providing code completions, inline documentation, editing suggestions, and simple refactorings. They are typically integrated directly into editors and IDEs, where developers interact with the underlying LLMs either through chat-like interfaces or by selecting code or comments using mouse-based interactions.

GitHub Copilot [89] and Cursor [10] are two representative examples of LLM-based code assistants. GitHub Copilot, originally developed based on Codex trained on GitHub code repositories, offers context-aware code completions across multiple programming languages and is tightly integrated into popular IDEs like Visual Studio Code [36] and JetBrains. Cursor extends this functionality by embedding conversational interaction, maintaining memory of previous edits, and supporting structured command execution. Other notable systems include Amazon Q Developer [44] and Tabnine [215]. Amazon Q Developer targets developers working in cloud ecosystems, offering language-specific completions, cloud API integration, and basic vulnerability detection. In contrast to Amazon Q Developer, Tabnine takes a privacy-first approach by deploying smaller local models trained on permissively licensed code, making it attractive for organizations and developers who do not want to send their code to a remote cloud.

Implementations of these systems typically exhibit reactive behavior, responding to user input without initiating their own plans or taking proactive steps. Their interactions are generally single-turn, relying on the immediate context within the code editor rather than maintaining a persistent memory of past interactions or broader development goals. Most systems in this category are tightly coupled with development tools and offer real-time assistance that fits naturally within existing programming workflows. However, they are limited in autonomy, lacking the ability to decompose complex tasks, maintain long-term state, or coordinate multi-step development activities.

Despite these limitations, interactive code assistants serve as a foundational layer within the more recent agentic programming ecosystem. They are widely deployed, easy to integrate into everyday development practices, and offer immediate value to developers. Furthermore, some recent systems, such as Cursor and GitHub Copilot Agent, are beginning to incorporate features like session-level memory, persistent context, and structured task execution, gradually bridging the gap between reactive code assistants and more autonomous, multi-turn agents.

These distinctions are captured in our taxonomy (see Table 5), where interactive code assistants are compared with other categories of agentic systems across key dimensions, including autonomy, memory scope, tool integration, reasoning complexity, and interaction model.

4.2.2 Autonomous task-oriented agents. These agents perform multi-step programming tasks with minimal human intervention, often maintaining control over the entire development process from requirement interpretation to code generation and validation. They can plan, execute, and revise their own workflows in response to intermediate results or changing task requirements. Many integrate external tools such as debuggers, package managers, or search engines, enabling them to gather information, resolve errors, and optimize code without direct user guidance. Examples include GPT-5 [181], which functions as an end-to-end collaborator for long-horizon code generation and debugging; Claude Opus 4 [51], engineered for sustained agentic workflows with strong long-context and tool-integration capabilities; Google Jules [18], which can sandbox repositories, propose diffs, and execute verified changes on real projects; DevGPT [75], an open-source pipeline that converts tickets into actionable code and CI artifacts; Kimi K2 [217], a model optimized for agentic coding proficiency; and Gemini 2 [82], which emphasizes multimodal inputs and built-in planning modes. These agents are often proactive in suggesting next steps, adaptive to new inputs, and capable of maintaining continuity across extended sessions through persistent memory or context management mechanisms.

4.2.3 Planning-centric agents. Planning-centric agents approach problem-solving as a two-phase process: first, structured task decomposition, where high-level goals are broken down into smaller, more manageable steps; and second, execution monitoring, in which results are evaluated and the plan is adjusted accordingly. This approach improves the handling of long-horizon tasks. For example, CAMEL [141] employs two agents in a role-playing setup, a “user” agent and an “assistant” agent, which collaboratively refine goals and strategies, producing plans that downstream code generators can execute. Voyager [225] demonstrates this paradigm by continuously exploring an open-ended environment, generating executable skills, and refining them into reusable plans for long-term adaptability. Similarly, CodePlan [232] introduces code-form planning, where structured pseudocode serves as an explicit intermediate representation to decompose complex problems into executable steps. These agents are typically multi-turn, memory-enabled, and trading speed for robustness.

4.2.4 Multi-agent and collaborative systems. Multi-agent and collaborative systems extend agent-based programming by introducing multiple specialized agents that coordinate to solve complex software engineering tasks. This approach draws inspiration from human software teams, where each member has a distinct role, such as requirements analysis, coding, testing, or documentation, and communication protocols are established to ensure progress toward shared objectives. For example, SWE-Agent [242] employs multiple role-specific LLM agents: an “Architect” agent for high-level design, a “Coder” agent for implementation, and a “Reviewer” agent for quality assurance, which are connected through structured dialogue and shared memory. ChatDev [191] follows this paradigm by simulating an end-to-end software company, where agents take on roles such as CEO, CTO, and programmers to collaboratively design, implement, and test applications. Furthermore, AutoCodeRover [250] extends collaboration into real-world repositories, orchestrating specialized agents that autonomously navigate, edit, and validate source code across complex multi-file projects.

4.3 Summary and Comparative Table

We conclude this section with a comparative summary (Table 5) of representative systems across the behavioural dimensions and categories introduced above. This taxonomy provides a lens to

Table 5. Comparison of representative AI agentic programming systems.

System	Category	Proactivity	Multi-turn	Tool Use	Adaptivity
GitHub Copilot [89]	IDE Assistant	Reactive	✗	✓	✗
Amazon Q Developer [44]	IDE Assistant	Reactive	✓	✓	✓
Tabnine [215]	IDE Assistant	Reactive	✗	✗	✗
Cursor [10]	IDE Assistant	Reactive	✗	✓	✗
Claude Opus 4-powered agent [51]	Task-oriented	Proactive	✓	✓	✓
Google Jules [18]	Task-oriented	Proactive	✓	✓	✓
DevGPT [75]	Task-oriented	Proactive	✓	✓	✓
KimI K2-powered agent [217]	Task-oriented	Proactive	✓	✓	✓
Gemini 2-powered agent [82]	Task-oriented	Proactive	✓	✗	✗
Voyager [225]	Planning Agent	Proactive	✓	✓	✓
CAMEL [141]	Planning Agent	Proactive	✓	✗	✗
CodePlan [232]	Planning Agent	Proactive	✓	✗	✗
ChatDev [191]	Multi-agent System	Proactive	✓	✓	✗
AutoCodeRover [250]	Multi-agent System	Proactive	✓	✓	✓
SWE-Agent [242]	Multi-agent System	Proactive	✓	✓	✓

Table 6. Example token pricing for LLMs used in coding tasks (USD per 1M tokens, as of Sep. 2025).

Model	Input(\$)	Output(\$)	Context Window
GPT-5 (Standard)	1.25	10.00	256k
GPT-5 Mini	0.25	2.00	128k
GPT-4 variants (e.g., 4o)	2.50	10.00	128k
Claude 4 Opus	15.00	75.00	200k
Gemini 2.5 Pro	1.25	10	1M
Grok 4	3.00	15.00	256k
DeepSeek R1-0528	0.55	2.19	160k
Kimi K2	0.15	2.50	128k
Qwen3-235B-A22B	0.22	0.88	128k
Qwen3-Coder-480B-A35B-Instruct	4.5	13.5	256k
Solar-Pro	0.30	0.30	128k
Openhands-LM-32B-v0.1	2.6	3.4	128k
Devstral-Small	0.10	0.30	128k
Devstral-Medium	0.40	2.00	128k

understand the capabilities and limitations of existing approaches and can guide the design of future systems.

4.4 Cost and Token Consumption Model

While recent LLMs show impressive capabilities in software engineering tasks, their real-world applicability is often constrained by cost considerations. This cost is typically measured in terms of *tokens consumed per US dollar* for both input and output, along with additional expenses incurred by extended reasoning strategies such as Chain-of-Thought (CoT) and tool-augmented workflows.

4.4.1 Pricing dimensions. Commercial providers generally price LLM usage by input and output tokens, with rates varying across model families. Some, such as OpenAI’s GPT-5, offer multiple service tiers (Standard, Mini, Nano, Pro) with different pricing, context lengths, and throughput limits. Table 6 summarizes representative pricing.

4.4.2 Impact of reasoning strategies. To make the cost-performance trade-offs more concrete, we draw on the “Agentic workflow for implementing a REST task” example described earlier in Section 2. In that workflow, the agent receives a high-level natural language specification for a REST API endpoint, interprets the requirements, generates code, and iteratively tests the implementation until it passes the provided unit tests.

We consider three reasoning strategies applied to this scenario:

Short reasoning. The model produces the implementation in a single turn with minimal intermediate reasoning. For the REST API task, this means directly generating the endpoint code and tests without explicit planning or validation steps. This approach minimizes token usage and latency but risks missing subtle requirements.

Standard CoT. The model uses a fixed-depth chain of thought to plan the implementation. In the REST API case, this involves reasoning about request handling, data validation, and error responses before generating the code. This strategy consumes significantly more tokens, roughly twice as many, compared to the short reasoning strategy, but yields a higher likelihood of producing a correct implementation on the first attempt.

Tool-augmented iterative reasoning. The agent integrates code compilation and test execution into the workflow. After producing an initial version, it runs the tests, inspects any failures, and revises the code in subsequent turns. For the REST API example, this may involve multiple cycles of fixing logic errors, adjusting request parsing, and refining edge-case handling. While this maximizes accuracy and robustness, it also increases token consumption and wall-clock time substantially due to repeated code generation and analysis.

In practice, the optimal choice depends on the cost-performance budget of the project. For time-sensitive or budget-constrained environments, a hybrid approach can offer a more effective balance.

5 Challenges

AI agentic programming introduces a promising and complex shift in how software is developed, relying on the autonomous capabilities of LLMs. Despite recent progress, several technical and conceptual challenges remain that hinder the deployment of robust, scalable, and trustworthy agentic systems [39, 199].

5.1 Evaluation and Benchmarking

A variety of benchmarks and open-source toolkits [153, 193] have been proposed to evaluate the capabilities of LLM-based agents across different tasks. Despite this progress, existing benchmarks commonly used for assessing coding agents, such as HumanEval [254] and SWE-Bench [121], may still be inadequate for capturing the full complexity of real-world software engineering workflows.

Most of these benchmarks primarily focus on small, self-contained problems, often restricted to a small number of programming languages, such as Python [62], and generally lack support for interactive, multi-turn, or tool-integrated tasks [153]. In contrast, practical agentic systems are expected to operate over large, modular codebases, interface with third-party libraries, manage build workflow pipelines, and respond dynamically to user feedback or runtime tool outputs. As LLM-based systems increasingly incorporate reinforcement learning [183, 226] and more advanced planning mechanisms, future benchmarks should reflect this integration. Table 7 shows the characteristics of tasks in SWE-Bench. Although SWE-Bench introduces project-level repositories and leverages unit tests and continuous integration (CI) for evaluation, its scope is restricted to Python. Most tasks are function- or module-level, with no support for multi-turn feedback, third-party library usage, or build pipeline management. Even for project-level tasks, interaction is limited to binary pass/fail outcomes, providing only minimal support for realistic multi-step or tool-integrated software engineering workflows.

Moreover, there is a noticeable absence of evaluation frameworks designed for emerging complex use cases, such as those involving interactions with compilers and debuggers [229], where agents must reason about low-level program behavior, perform iterative transformations, or track state

Table 7. SWE-Bench characteristics.

Task type	Proportion	Languages	Interaction	Multi-turn feedback	Library integration	Build pipeline
Function-level	65%	Python	Unit tests	✗	✗	✗
Module-level	25%	Python	Tests and CI	✗	✗	✗
Project-level	<10%	Python	Tests and CI	Limited (pass/fail)	✗	✗

across toolchains. The lack of such domain-specific benchmarks presents a significant gap in evaluating agent performance under realistic conditions.

5.2 Communication Protocols for Multi-Agent Systems

Early protocols enabled one-to-one agent-to-tool interactions but did not support direct agent-to-agent communication [25]. Current practice often relies on heterogeneous web service protocols and adapters, which provide interoperability but introduce high latency, bandwidth overhead, and limited scalability [172]. Recent work [1, 224, 239] explores group session models, where each execution unit functions as a session or group agent coordinating multiple services, and standards efforts propose unified message formats and session semantics. However, the lack of a common protocol, inefficiencies of adapters, and challenges in managing dynamic multi-agent sessions remain open problems for scalable and dependable agentic systems.

5.3 Domain-specific Models for Agents

Generic coding agents often struggle in domain-specific environments such as embedded systems, high-performance computing, optimization, or formal software verification [137, 159]. These domains typically impose stricter operational constraints and require deep integration with specialized APIs, toolchains, and domain knowledge resources that are often underrepresented in general-purpose training corpora. To address these limitations, recent research has proposed domain-adapted models and task-specific learning strategies to accelerate agent performance in specialized settings [208]. For instance, some approaches have begun incorporating compiler knowledge or security-specific patterns into LLM training pipelines [69, 148], enabling agents to reason more effectively about low-level program behavior or vulnerability patterns. In the future, developing robust and adaptable domain foundation models will be a promising direction for enabling agents to operate reliably in complex software environments, such as LLMs pretrained or fine-tuned on domain-specific data, tools, and semantics.

5.4 Safety and Privacy

As agentic systems gain increasing autonomy, so does the potential for unsafe behavior. Unlike traditional tools [6], agentic systems can invoke external tools, perform structural code modifications, and even commit changes without direct human oversight [221, 251]. These capabilities introduce significant risks, including the possibility of introducing subtle bugs, propagating unsafe patterns, or violating security constraints. A critical future direction involves ensuring that agentic systems can protect users and data. For example, when agents visit private repositories or are deployed in cloud-integrated environments, future models may need built-in controls to restrict access to sensitive project data. Further, malicious prompts, poisoned APIs, or compromised toolchains can mislead agents into executing unsafe behaviors [72, 247]. Future research should prioritize the design of secure protocols for agent collaboration, including authentication between agents, validation of tool outputs, and detection of anomalous actions. Also, agents must be capable of explaining their reasoning, flagging uncertainties, and allowing developers to understand and

revise with minimal effort. Building safety and privacy into the foundation of agentic architectures is essential.

5.5 Toolchain Integration and Programming Language Design

One fundamental challenge is the incompatibility between existing software tools and the needs of autonomous coding agents. Most programming languages, compilers, debuggers, and development environments were designed for human developers [198]. They emphasize usability and readability over structured, machine-consumable feedback. As a result, agents often struggle to diagnose failures, trace the consequences of code transformations, or interpret build errors [127]. For example, compilers typically report transformation failures or semantic conflicts in the form of coarse error messages, providing little insight into why an optimization was blocked or how a type error emerged [143]. Languages similarly prioritize human readability over machine-negotiated meaning, while compilers conceal internal reasoning to avoid overwhelming human users. These design choices, while historically effective, now limit the ability of AI coding agents to construct safe, efficient, and adaptive software.

To enable agentic workflows, toolchains should evolve to expose richer intermediate representations, transformation traces, and structured feedback interfaces. Equally important, programming languages should incorporate annotations and agent-aware interfaces that make developer intent explicit, allowing automated reasoning to be guided by semantic contracts rather than inferred heuristics. Together, these advances would transform compilers and languages from opaque tools into collaborative infrastructures capable of supporting autonomous agents at scale.

5.6 Scalable Memory

Agentic programming systems must maintain coherence and reasoning over long-running tasks involving multiple iterations, tools, and contextual dependencies. However, current LLMs are limited by fixed context windows and lack persistent, structured memory mechanisms [228]. Realistic software tasks may require agents to store and reason over evolving states, feedback logs, intermediate plans, and prior actions [192]. Without hierarchical and queryable memory systems, agents risk repeating errors, forgetting past successes, or producing inconsistent results. Emerging solutions such as retrieval-augmented generation and memory summarization offer partial relief, but they remain inadequate for complex, multi-session workflows [162]. Future research can explore memory architectures that differentiate short-term interactions, mid-term subgoals, and long-term domain knowledge.

6 Opportunities and Future Directions

AI agentic programming represents a fast-evolving research frontier that intersects artificial intelligence, programming languages, and software engineering. While recent advances have demonstrated promising capabilities, significant challenges remain in realizing robust, efficient, and trustworthy agentic systems. In this section, we outline several key opportunities and open research directions that can shape the future of this field. An overview is provided in Figure 8.

6.1 Integrating Coding Agents with Tools

Existing AI coding agents typically orchestrate LLMs with loosely integrated toolchains and basic memory mechanisms. These ad hoc designs often lack robustness, scalability, and generalization across programming tasks. Advancing agent architectures will require moving beyond simple prompt-response patterns toward more modular, structured systems that support reasoning, tool interaction, planning, and verification.

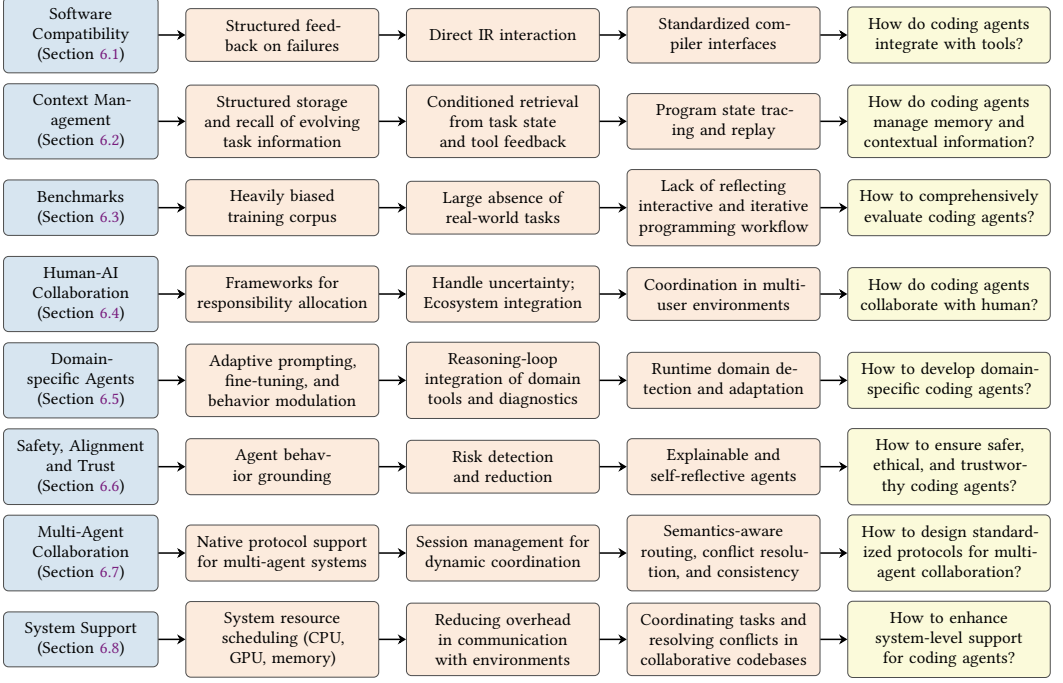


Fig. 8. Summary of potential future research directions of AI agentic programming.

A promising direction is to rethink how programming languages, compilers, and testing frameworks, which are traditionally built for human developers, can be redesigned to support AI coding agents. For example, instead of emitting opaque diagnostics, compilers could provide structured feedback explaining why certain optimizations (e.g., vectorization or inlining) fail [53, 163, 206]. These could include semantic barriers like unresolved aliasing, ambiguous data/control flow, or missing annotations [55], enabling agents to revise code more precisely. Beyond diagnostics, compilers can help agents track state across iterations. Feedback on which edits introduced errors, failed assertions, or performance regressions would enable agents to reason over the change history and adjust strategies accordingly.

Opening compiler internals also presents a valuable opportunity. Coding agents could interact directly with intermediate representations (IRs), such as LLVM IR [134] or MLIR [135], to reason about program structure, verify transformations, or perform static analysis at a semantic level. Compiler APIs and language servers (e.g., Clang’s LibTooling, the Language Server Protocol) already expose ASTs, symbol tables, and refactoring tools, but a wider adoption may require standardized, introspective interfaces across compilers.

At the programming language level, agent-aware extensions or annotations could further improve interaction. Developers might use domain-specific languages, embedded contracts, or even natural language comments, e.g., “sort the elements of input array x ”, to convey intent. This could guide synthesis, verification, or debugging. Likewise, compilers might expose symbolic summaries of control flow, memory access patterns, or performance profiles to inform multi-step agent reasoning.

Tighter integration with runtime systems also offers opportunities. For instance, agents can dynamically insert instrumentation or launch profiling runs, then use the results to inform optimization choices. Coupling these capabilities with autotuning frameworks [47, 63] would expand

the design space while preserving correctness and safety. Finally, advances in structured code representations, such as ASTs, graph-based IRs, and semantic embeddings, offer a foundation for more powerful agent reasoning. Combining LLMs with graph neural networks or neuro-symbolic systems could improve generalization and support cross-language, cross-target understanding.

6.2 Scalable Memory and Context Management

A key capability of agentic programming lies in managing memory and contextual information across tasks that involve long context reasoning and multiple iterations. Unlike traditional code generation, which typically follows a single pass prompt to solution model, agentic workflows for solving real-world software engineering problems involve multiple steps, iterative refinement, and integration with external tools and development environments [10, 89, 179, 230].

Consider an agent tasked with adding a new feature to an open source project, such as implementing a command-line flag to enable verbose logging. The agent must first analyze the existing codebase to locate the argument parsing logic, generate the required code changes, and update the logging behavior. If the updated code fails with a runtime error due to an uninitialized flag, the agent needs to debug the issue by inspecting stack traces, revise the code accordingly, and rerun the tests. Once the implementation is verified, the agent writes a commit message, creates a pull request with a summary of the changes, and links it to the relevant issue.

Throughout this process, the agent must persist and reason over a large and evolving context: the initial task description, previously generated code, compiler and runtime feedback, and version control metadata. Without the ability to store and recall this information in a structured way, the agent may repeat past mistakes, forget earlier successful changes, or submit incomplete solutions. Agentic programming, therefore, needs mechanisms for memory and context tracking that go beyond simple token limits, enabling coding agents to maintain continuity across extended interactions and tool usage.

As the memory footprint of LLMs grows linearly with input token length [126, 223], current LLM-based agents remain constrained by their context windows and lack persistent memory across a long sequence of iterations [152, 184, 200]. However, reasoning about real-world programs often requires modeling complex data structures and code context (like function calls) spanning across multiple files, which often exceeds typical context limits. Although some industry-scale LLMs claim to support million-token contexts [77, 78, 216], they often rely on random sampling techniques [102] and fail to leverage program structure or semantics effectively.

Therefore, an interesting direction is to design attention mechanisms that are guided by code structure, such as syntax trees, control flow graphs, or data dependencies. These structures can help agents focus more accurately on the most relevant parts of a program. While approaches like retrieval-augmented generation (RAG) [140], KV cache offloading [138, 211], and compression [60, 158] provide partial solutions, they struggle to provide precise control over long-term dependencies, structured knowledge, and execution histories.

AI coding agents can also benefit from hierarchical memory models that distinguish between short-term interaction history, mid-term planning objectives (such as subgoals and intermediate decisions), and long-term knowledge. This long-term layer may include patterns of success or failure, reusable code templates, and observed tool behaviors. Such hierarchies can be dynamically updated and selectively queried using retrieval controllers or attention-based mechanisms. Additionally, memory summarization techniques could be explored to condense lengthy interaction histories into structured, semantically meaningful representations. For example, an agent might summarize a multi-turn session as a sequence of planning decisions and outcomes, highlighting key insights and interventions.

Table 8. Benchmarks for evaluating LLMs and agentic systems on programming tasks. Abbreviations: CP = Competitive Programming, BF = Bug Fixing, FC = Function Completion, CR = CLI Reasoning, PO = Performance Optimization.

Benchmark	Source	Language	Task	Difficulty	Year
HumanEval	Hand-written	Python	FC	Beginner	2021
MBPP	Crowd-sourced	Python	FC	Beginner	2021
CodeContests	CP	Python/C++/Java	FC	Diverse	2022
HumanEval-X	Hand-written	Python/C++/Java/JS/Go	FC	Intermediate	2023
SWE-Bench	GitHub Issues	Python	BF	Expert	2024
SWE-bench M	GitHub Issues	JS	BF	Diverse	2024
LiveCodeBench	CP (live)	Python	FC, BF	Diverse	2024
Terminal-Bench	Community-curated	Shell	CR	Diverse	2025
Spider 2.0	Enterprise DB Apps	SQL	FC	Expert	2025
EffiBench-X	Synthetic	Python/C++/JS	FC, BF	Diverse	2025
Web-Bench	Web App Projects	JS/TS/HTML/CSS/Python	BF, FC	Expert	2025
ProjectEval	Open-source Repos	Python/Java/C++/JS	FC, BF, CR	Expert	2025
TRAIL	CP	Python/Java/C++/JS	BF, CR	Expert	2025
GSO	Open-source Repos	Python/C/C++/Cython/Rust	PO	Expert	2025

Note: "Diverse" difficulty indicates that the benchmark covers a wide range from beginner to expert-level tasks.

Another important area is the development of context-aware retrieval strategies that move beyond static similarity-based methods [209]. During debugging, for instance, an agent could retrieve not only the most recent error message but also similar past failures, proposed fixes, relevant test cases, and their outcomes. Retrieval conditioned on task state and tool feedback would significantly improve the agent’s ability to reason under uncertainty.

Structured mechanisms for program state tracing and replay may also enhance agent performance. By recording partial program states, tool outputs, and execution steps, agents can support backtracking, recovery from failure, and richer explanations. For example, an agent could explain how a specific code edit introduced a type error or why a particular memory access blocked loop vectorization. These capabilities are crucial for supporting causal reasoning and improving transparency. Likewise, persistent memory across multiple code generation and refinement sessions will be essential for enabling agents to accumulate and refine knowledge over time. This may include long-term storage of project-specific context, interaction histories, usage patterns of tools, and models of user intent. Such memory infrastructure will allow for continual learning and increasing personalization of agent behavior.

In summary, effective memory and context management are foundational for scaling agentic programming systems. These capabilities are vital for advancing from reactive, prompt-driven assistants to autonomous, context-aware collaborators capable of sustained reasoning, adaptation, and long-term learning.

6.3 Evaluation and Benchmarking

Table 8 summarizes several widely used coding benchmarks for evaluating LLMs and agentic systems on programming tasks. These benchmarks vary in their origin, programming language coverage, task types, and difficulty levels. Commonly used datasets include *HumanEval* [179], *HumanEval-X* [254], *MBPP* [54], *SWE-Bench* [121], *SWE-bench Multimodal* [242], *Terminal-Bench* [218], *LiveCodeBench* [115], *CodeContests* [146], *Spider2.0* [139], *EffiBench-X* [194], *Web-Bench* [240], *ProjectEval* [151], *TRAIL* [74], and *GSO* [204].

While these benchmarks have provided valuable insights into the capabilities of LLMs and agentic systems for code generation and bug fixing, they also have important limitations. For example, they are heavily biased toward a small set of programming languages - especially Python, which dominates the training corpus of code for current models [59, 222]. This limits the generalizability of evaluation results to domains involving statically typed or domain-specific languages, such

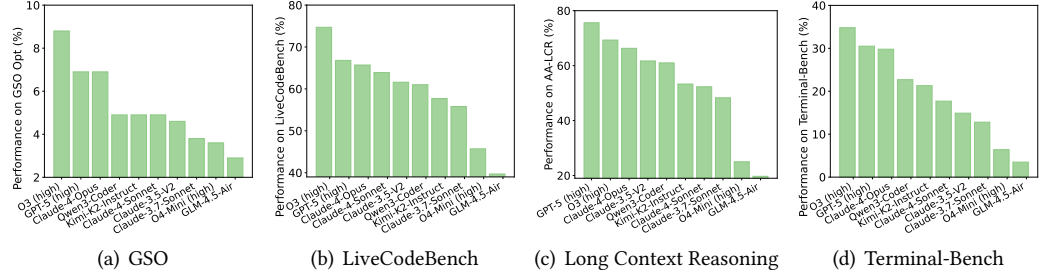


Fig. 9. Performance comparison of the top 10 language models across four representative benchmarks, including software optimization (GSO), code generation (LiveCodeBench), long-context reasoning (AA-LCR), and system-level tasks (Terminal-Bench).

as C++ or Rust [222]. For example, Figure 9 presents the performance of the top-10 LLMs across four representative benchmarks. The results reveal clear task-dependent differences. For example, all models still perform poorly on software optimization tasks, while they perform well on code generation tasks. Overall, performance gaps remain across all benchmarks, indicating that current LLMs still face significant challenges in reasoning, coding, and generalization to complex tasks.

Furthermore, many existing benchmarks focus on small, self-contained problems that may not be representative of real-world software engineering tasks. Realistic development scenarios often involve working with large, modular codebases, extensive use of third-party libraries, non-trivial build processes, and long-range dependencies across files and components. These characteristics are largely absent from current benchmark datasets, making it difficult to assess an agent’s ability to scale or generalize. Another key limitation is that most benchmarks do not capture the interactive, iterative nature of agentic programming. In real-world settings, coding agents will need to collaborate with human developers [10, 89, 179, 230], receive intermediate feedback and confirmation, and rely on external tools such as compilers, debuggers, and test frameworks. Benchmarks that assume single-shot or non-interactive task completion fail to reflect the complexity of such multi-step, tool-augmented workflows.

Addressing these gaps will require developing more comprehensive and extensible evaluation frameworks. Future benchmarks should incorporate realistic tasks that reflect end-to-end development workflows, support multiple programming languages, and enable interaction with tools and human feedback loops. Metrics should go beyond functional correctness to include robustness, tool usage efficiency, recovery from failure, and the ability to incorporate feedback. Simulation environments and evaluation harnesses will also be important for reproducibility and fair comparison.

6.4 Human-AI Collaboration

While a long-term vision of agentic programming is to automate the entire software development lifecycle, including writing, debugging, and testing code, near-term opportunities include extending the capabilities of current LLM-based coding assistants [10, 89, 179, 230]. Rather than replacing human developers, these systems can act as collaborative partners, supporting workflows in which humans retain strategic oversight. Similar to pair programming [97], LLM-based agents can assist by proposing ideas, detecting errors, suggesting improvements, and automating routine tasks, thereby augmenting human productivity.

Designing effective models for human-AI collaboration is a key research challenge. This involves creating user interfaces, interaction protocols, and responsibility-sharing frameworks suited for professional and team-based environments where coordination, trust, and efficiency are critical.

Unlike traditional code generation tools, agentic systems operate in iterative, feedback-driven loops and make autonomous decisions based on intermediate results. This shifts the developer's role from command giver to collaborator, opening new opportunities and challenges for co-creating software. To support this shift, agents must be transparent in their reasoning, responsive to human input, and capable of explaining their decisions, which calls for advances in interactive prompting [201], natural language explanations [68], and context-aware dialogue protocols [48, 203].

An interesting direction would be mixed initiative workflows, where control moves fluidly between human and agent. Developers can specify goals, constraints, or coding conventions, while agents generate scaffolding code, explore alternatives, or automate repetitive tasks. The developer then inspects, accepts, rejects, or refines the output. For example, an agent may suggest several refactoring options and generate test cases to validate them, leaving the final choice to the human. Such workflows must also support interruption and resumption, enabling agents to incorporate new input without losing context. This requires toolchains that track shared state across iterations, including goals, partial programs, and transformation history.

Finally, domain-aware collaboration can strengthen interpretability and safety. Agents should handle ambiguous or incomplete requirements by asking clarifying questions [171, 248], drawing on design documents, issue trackers, and domain knowledge to refine their reasoning. They should also provide clear explanations, cite relevant evidence, and highlight trade-offs to build trust, which is especially critical in domains such as finance, healthcare, and aerospace.

6.5 Domain Specialization and Adaptability

While LLMs trained on extensive code corpora can demonstrate general programming capabilities [122, 253], they may struggle in specialized domains, such as generating Verilog code for hardware synthesis, which involves strict performance constraints, low-level abstractions, or domain-specific libraries and tooling.

Future research can improve agentic programming by exploring domain adaptive prompting, fine-tuning, and behavior modulation in areas such as embedded systems, data science, scientific computing, high-performance computing, and formal methods. In embedded systems programming [81], for example, agents must reason about memory layout, timing constraints, and hardware-specific APIs, which are often underrepresented in general training data. Training on domain-specific code or adjusting planning strategies to prioritise correctness and safety over brevity may enhance performance. Integrating domain-specific tools and diagnostics into agents' reasoning, such as simulation engines [182, 189, 219], static analyzers [132], fuzzers [165], or formal verification tools [234], can further improve robustness. Specialisation also supports interpretability and safety, enabling stricter validation, better use of formal specifications, and more reliable feedback in safety-critical domains such as aerospace, automotive, and healthcare.

Additionally, research into adaptive agent behavior can enable models to detect the domain of a task at runtime and adjust their prompting strategies, tool usage, and explanation styles accordingly. For example, an agent working on a financial modeling script might switch to using Pandas and SQL queries, while one dealing with real-time control code may emphasize low-latency function design and interrupt safety. This line of research will also open opportunities for personalisation at the developer or team level. Agents can learn preferences, coding conventions, and domain-specific heuristics from local repositories or past interactions, enabling more consistent and context-aware assistance.

6.6 Safety, Alignment, and Trust

As agentic coding systems start taking on greater responsibility in software development, it becomes increasingly important to ensure that their behavior aligns with user intent, produces correct results,

and avoids unintended changes [149, 183]. Unlike traditional LLM-based code assistants, agentic systems can take multi-step actions, use external tools, and modify codebases with limited human oversight, making safety and trust essential goals [104, 178].

One key research direction is building agents that can better understand and follow what users actually want, even when instructions are vague or incomplete. Current systems rely heavily on natural language prompts, which can be ambiguous. Future work could focus on grounding agent behavior in more structured forms of input, such as constraints, test cases, or high-level goals, that are easier to validate and reason about [83, 107, 142]. Another interesting idea is to design a structured language that developers can use to express their intent clearly. This language would serve as a kind of programming interface for LLMs, allowing users to define what the agent should do, what it must avoid, and what counts as a valid solution. For example, a developer might specify that a function's output must remain the same after refactoring, and the agent would only explore changes that meet this requirement. This kind of structure could also make it easier to apply lightweight verification tools, such as static analyzers or type checkers, to ensure the agent's suggestions are safe and correct [132, 147, 177].

In safety-critical domains like healthcare, finance, or automotive software, agents will also need to follow strict coding standards, legal rules, and certification guidelines. Agents could be trained to recognize such constraints or be paired with rule-based systems to flag violations and suggest compliant alternatives [98, 183]. Another important challenge is ethical alignment. Because agentic systems are trained on large and diverse codebases, they may learn unsafe or biased practices [106, 154, 187]. Research efforts are needed to detect and reduce risks such as generating insecure code, leaking sensitive data, or reinforcing stereotypes. Techniques from responsible AI, like behavioral audits, adversarial testing, and human feedback, can be adapted to this setting.

To build user trust, agentic programming systems should also be able to explain their decisions [68, 196, 205]. Developers want to know why an agent made a change, where the idea came from, and what trade-offs are involved. Research on explanation generation, source citation, and visualization of code changes will help make LLM behavior more transparent and understandable [108, 156, 213]. Agents should also be aware of their limitations. When they are uncertain, they should be able to flag their confidence level [68, 130], suggest multiple options [161, 205], or ask the user for confirmation [171, 248]. Designing agents that can adapt their behavior based on task difficulty or user feedback is a promising direction [61, 71, 166, 205, 212]. Furthermore, safety mechanisms like undo and audit trails will be essential. Agents should keep track of what they have changed and allow users to roll back actions easily [120, 186]. Future work could explore automatic snapshotting, reversible code edits, and tight integration with version control systems to support safe collaboration.

6.7 Multi-agent Collaboration

As coding agents increasingly operate in teams and interact with heterogeneous services, effective multi-agent collaboration requires standardized communication protocols that provide low-latency, scalable, and semantics-aware coordination. However, current solutions remain limited and fall short of these needs. Future communication protocols should move beyond ad-hoc adapters and heterogeneous web services toward native support for multi-agent collaboration. Promising directions include the design of lightweight, low-latency group communication frameworks that natively support both agent-to-agent and agent-to-tool interactions, while offering standardized session management for dynamic participation and coordination [1, 224, 239]. Recent standards efforts such as IEEE P3394 propose a universal framework with session, channel, and message abstractions, enabling consistent semantics and interoperability across diverse transports [220]. Building on these initiatives, future research should emphasize semantics-aware routing, conflict

resolution, and consistency guarantees in distributed multi-agent settings. Ultimately, advancing collaboration-centric protocols will be essential for scaling efficient agentic systems in practice.

6.8 System Support for AI Coding Agents

Realizing agentic programming at scale also requires strong systems-level support. When offered as a cloud service [43, 110, 129, 169], LLMs and coding agents must be backed by infrastructures that can efficiently manage resources, communication, and state. Agentic workflows combine code generation, compilation, testing, and tool orchestration, with highly variable resource demands [230, 242]. Systems must schedule CPU, GPU, and memory adaptively, balance throughput across users, and minimize latency for interactive tasks [45, 87, 185, 244, 245]. During the iterative code generation and testing process, agents are likely to exchange large volumes of data with tools like compilers, debuggers, profilers, and repositories. Reducing overhead through caching [131, 155], compression [101, 118, 241, 252], and lightweight protocols [93, 170] is essential for efficiency in distributed environments. Equally important is context and state management. To support interruption, resumption, and collaboration, systems must persist not only source code but also goals, annotations, and intermediate representations. Efficient checkpointing and synchronization are key enablers here. Finally, multi-agent collaboration introduces distributed systems challenges: coordinating tasks, resolving conflicts, and ensuring consistency when multiple agents - or human-agent teams - work concurrently on shared codebases.

7 Conclusion

We have presented a comprehensive review of AI agentic programming. This new software paradigm, driven by the recent success of LLMs, is a transformative shift in how software can be created, maintained, and evolved. By combining the capabilities of LLMs with planning, tool use, and iterative refinement, coding agents are beginning to automate complex, multi-stage programming workflows that traditionally required significant human involvement. These systems can not only generate and test code but also interact with development tools, decompose tasks, and adapt based on feedback, bringing us closer to the vision of autonomous software development.

We have introduced a taxonomy of AI coding agents and architectures, reviewed underpinning techniques like context management and tool integration, and discussed how existing benchmarks evaluate the capabilities of coding agents. We have summarized the progress made in enabling LLM-based AI agents to reason over tasks, interface with software development tools, and operate in increasingly sophisticated ways. At the same time, we have identified multiple open challenges that must be addressed to ensure these systems are safe, reliable, and usable in real-world settings. These include limitations in context handling, the need for persistent and structured memory, alignment with user intent, human-AI collaboration, and verification of agent behavior. As software developers are increasingly relying on AI coding agents, these concerns will become more pressing and will require interdisciplinary solutions drawing from programming languages, human-computer interaction, software engineering, and responsible AI.

Looking ahead, AI agentic programming offers exciting opportunities to fundamentally rethink the software development practice. Whether as collaborative partners in interactive workflows or as autonomous systems that manage long-running tasks, these agents have the potential to augment developer productivity, reduce software maintenance costs, and expand access to programming. We hope this survey serves as a foundation for researchers and practitioners to navigate the emerging landscape of AI agentic programming and to accelerate progress in building the next generation of intelligent and trustworthy software development tools.

References

- [1] 2025. Agent-to-Agent (A2A) Protocol. <https://github.com/a2aproject/A2A>. Accessed: 2025-09-15.
- [2] 2025. Apache Maven: Project Management Tool. <https://maven.apache.org/>. Accessed: 2025-09-15.
- [3] 2025. Black: The uncompromising Python code formatter. <https://black.readthedocs.io/>. Accessed: 2025-09-15.
- [4] 2025. Claude Opus 4. <https://www.anthropic.com/claude/opus>. Accessed: 2025-09-15.
- [5] 2025. Cargo: Rust's package manager and build system. <https://doc.rust-lang.org/cargo/>. Accessed: 2025-09-15.
- [6] 2025. Checkmarx. <https://checkmarx.com/>.
- [7] 2025. Clang: a C language family frontend for LLVM. <https://clang.llvm.org/>. Accessed: 2025-09-15.
- [8] 2025. CMake: Cross-Platform Make. <https://cmake.org/>. Accessed: 2025-09-15.
- [9] 2025. Continue. <https://www.continue.dev/>. Accessed: 2025-09-15.
- [10] 2025. Cursor. <https://cursor.com>. Accessed: 2025-09-15.
- [11] 2025. Devika. <https://github.com/stitionai/devika>. Accessed: 2025-09-15.
- [12] 2025. ESLint: Find and fix problems in your JavaScript code. <https://eslint.org/>. Accessed: 2025-09-15.
- [13] 2025. flake8: the modular source code checker. <https://flake8.pycqa.org/>. Accessed: 2025-09-15.
- [14] 2025. GDB: The GNU Project Debugger. <https://www.gnu.org/software/gdb/>. Accessed: 2025-09-15.
- [15] 2025. Git: Distributed version control system. <https://git-scm.com/>. Accessed: 2025-09-15.
- [16] 2025. GNU Compiler Collection (GCC). <https://gcc.gnu.org/>. Accessed: 2025-09-15.
- [17] 2025. GNU Make. <https://www.gnu.org/software/make/>. Accessed: 2025-09-15.
- [18] 2025. Google Jules. <https://jules.google/>. Accessed: 2025-09-15.
- [19] 2025. IntelliSense. <https://code.visualstudio.com/docs/editing/intellisense>. Accessed: 2025-09-15.
- [20] 2025. javac: The Java Compiler. <https://docs.oracle.com/javase/8/docs/technotes/tools/windows/javac.html>. Accessed: 2025-09-15.
- [21] 2025. Jest: Delightful JavaScript Testing. <https://jestjs.io/>. Accessed: 2025-09-15.
- [22] 2025. Language Server Protocol. <https://microsoft.github.io/language-server-protocol/>. Accessed: 2025-09-15.
- [23] 2025. LLDB Debugger. <https://lldb.llvm.org/>. Accessed: 2025-09-15.
- [24] 2025. Mocha: JavaScript test framework. <https://mochajs.org/>. Accessed: 2025-09-15.
- [25] 2025. Model Context Protocol (MCP). <https://docs.anthropic.com/en/docs/mcp>. Accessed: 2025-09-15.
- [26] 2025. npm: Node Package Manager. <https://www.npmjs.com/>. Accessed: 2025-09-15.
- [27] 2025. OpenDevin. <https://github.com/AI-App/OpenDevin.OpenDevin>. Accessed: 2025-09-15.
- [28] 2025. pdb — The Python Debugger. <https://docs.python.org/3/library/pdb.html>. Accessed: 2025-09-15.
- [29] 2025. pip: The Python Package Installer. <https://pip.pypa.io/>. Accessed: 2025-09-15.
- [30] 2025. Prettier: An opinionated code formatter. <https://prettier.io/>. Accessed: 2025-09-15.
- [31] 2025. Pyright: Fast type checker for Python. <https://github.com/microsoft/pyright>. Accessed: 2025-09-15.
- [32] 2025. pytest: simple powerful testing with Python. <https://docs.pytest.org/>. Accessed: 2025-09-15.
- [33] 2025. tsserver: TypeScript language server. [https://github.com/microsoft/TypeScript/wiki/Standalone-Server-\(tsserver\)](https://github.com/microsoft/TypeScript/wiki/Standalone-Server-(tsserver)). Accessed: 2025-09-15.
- [34] 2025. TypeScript Compiler (tsc). <https://www.typescriptlang.org/docs/handbook/compiler-options.html>. Accessed: 2025-09-15.
- [35] 2025. unittest — Unit testing framework (Python). <https://docs.python.org/3/library/unittest.html>. Accessed: 2025-09-15.
- [36] 2025. Visual Studio Code. <https://code.visualstudio.com/>. Accessed: 2025-09-15.
- [37] 2025. Yarn: Fast, reliable, and secure dependency management. <https://yarnpkg.com/>. Accessed: 2025-09-15.
- [38] Emmanuel Abbe, Samy Bengio, Aryo Lotfi, Colin Sandon, and Omid Saremi. 2024. How far can transformers reason? the globality barrier and inductive scratchpad. *Advances in Neural Information Processing Systems* 37 (2024), 27850–27895.
- [39] Deepak Bhaskar Acharya, Karthigeyan Kuppan, and B Divya. 2025. Agentic ai: Autonomous intelligence for complex goals—a comprehensive survey. *IEEE Access* (2025).
- [40] Alfred V. Aho, Monica S. Lam, Ravi Sethi, and Jeffrey D. Ullman. 2006. *Compilers: Principles, Techniques, and Tools* (2nd Edition). Addison-Wesley Longman Publishing Co., Inc., USA.
- [41] Miltiadis Allamanis, Earl T. Barr, Premkumar Devanbu, and Charles Sutton. 2018. A Survey of Machine Learning for Big Code and Naturalness. *ACM Comput. Surv.* 51, 4, Article 81 (July 2018), 37 pages. doi:10.1145/3212695
- [42] Nadia Alshahwan, Jubin Chheda, Anastasia Finogenova, Beliz Gokkaya, Mark Harman, Inna Harper, Alexandru Marginean, Shubho Sengupta, and Eddy Wang. 2024. Automated Unit Test Improvement using Large Language Models at Meta (FSE 2024). Association for Computing Machinery, New York, NY, USA, 185–196. doi:10.1145/3663529.3663839
- [43] Amazon. 2025. Amazon SageMaker. <https://aws.amazon.com/sagemaker/>. Accessed: 2025-09-11.
- [44] Amazon Web Services. 2025. Amazon Q Developer: The most capable generative AI-powered assistant for software development. <https://aws.amazon.com/q/developer/>. Accessed: 2025-09-15.

- [45] Reza Yazdani Aminabadi, Samyam Rajbhandari, Ammar Ahmad Awan, Cheng Li, Du Li, Elton Zheng, Olatunji Ruwase, Shaden Smith, Minjia Zhang, Jeff Rasley, et al. 2022. Deepspeed-inference: enabling efficient inference of transformer models at unprecedented scale. In *SC22: International Conference for High Performance Computing, Networking, Storage and Analysis*. IEEE, 1–15.
- [46] Rohan Anil, Sebastian Borgeaud, and et al. 2023. Gemini: A Family of Highly Capable Multimodal Models. (2023). arXiv:2312.11805
- [47] Jason Ansel, Shoaib Kamil, Kalyan Veeramachaneni, Jonathan Ragan-Kelley, Jeffrey Bosboom, Una-May O'Reilly, and Saman Amarasinghe. 2014. Opentuner: An extensible framework for program autotuning. In *Proceedings of the 23rd international conference on Parallel architectures and compilation*. 303–316.
- [48] Anthropic. 2024. Introducing the Model Context Protocol (MCP). <https://www.anthropic.com/news/model-context-protocol>. Accessed: 2025-09-15.
- [49] Anthropic. 2024. Introducing the next generation of Claude. <https://www.anthropic.com/news/claude-3-family>. Accessed: 2025-09-15.
- [50] Anthropic. 2025. Claude Code: Command-Line Interface Reference. <https://docs.anthropic.com/en/docs/claude-code/cli-reference>. Accessed: 2025-09-15.
- [51] Anthropic. 2025. Claude Opus 4. <https://www.anthropic.com/claude/opus>. Accessed: 2025-09-15.
- [52] Anthropic. 2025. Model Context Protocol (MCP) Specification. <https://modelcontextprotocol.io/>. Accessed: 2025-09-15.
- [53] Amir H Ashouri, William Killian, John Cavazos, Gianluca Palermo, and Cristina Silvano. 2018. A Survey on Compiler Autotuning using Machine Learning. *Computing Surveys (CSUR)* 51 (2018), 1–42.
- [54] Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. 2021. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732* (2021).
- [55] Zhangqian Bi, Yao Wan, Zheng Wang, Hongyu Zhang, Batu Guan, Fangxin Lu, Zili Zhang, Yulei Sui, Hai Jin, and Xuanhua Shi. 2024. Iterative Refinement of Project-Level Code Context for Precise Code Generation with Compiler Feedback. In *Findings of the Association for Computational Linguistics ACL 2024*. 2336–2353.
- [56] Islem Bouzenia, Premkumar Devanbu, and Michael Pradel. 2024. Repairagent: An autonomous, llm-based agent for program repair. *arXiv preprint arXiv:2403.17134* (2024).
- [57] Marcel Bruch, Martin Monperrus, and Mira Mezini. 2009. Learning from examples to improve code completion systems. In *Proceedings of the 7th Joint Meeting of the European Software Engineering Conference and the ACM SIGSOFT Symposium on The Foundations of Software Engineering* (Amsterdam, The Netherlands) (*ESEC/FSE '09*). Association for Computing Machinery, New York, NY, USA, 213–222. doi:10.1145/1595696.1595728
- [58] Marcel Bruch, Martin Monperrus, and Mira Mezini. 2009. Learning from examples to improve code completion systems. In *Proceedings of the 7th joint meeting of the European software engineering conference and the ACM SIGSOFT symposium on the foundations of software engineering*. 213–222.
- [59] Linzheng Chai, Shukai Liu, Jian Yang, Yuwei Yin, Ke Jin, Jiaheng Liu, Tao Sun, Ge Zhang, Changyu Ren, Hongcheng Guo, et al. 2024. Mceval: Massively multilingual code evaluation. *arXiv preprint arXiv:2406.07436* (2024).
- [60] Chi-Chih Chang, Wei-Cheng Lin, Chien-Yu Lin, Chong-Yan Chen, Yu-Fang Hu, Pei-Shuo Wang, Ning-Chi Huang, Luis Ceze, Mohamed S Abdelfattah, and Kai-Chiang Wu. 2025. Palu: KV-Cache Compression with Low-Rank Projection. In *The Thirteenth International Conference on Learning Representations*.
- [61] Justin Chih-Yao Chen, Archiki Prasad, Swarnadeep Saha, Elias Stengel-Eskin, and Mohit Bansal. 2024. Magicore: Multi-agent, iterative, coarse-to-fine refinement for reasoning. *arXiv preprint arXiv:2409.12147* (2024).
- [62] Mark Chen, Jerry Tworek, et al. 2021. Evaluating Large Language Models Trained on Code. arXiv:2107.03374 [cs.LG] <https://arxiv.org/abs/2107.03374>
- [63] Tianqi Chen, Thierry Moreau, Ziheng Jiang, Lianmin Zheng, Eddie Yan, Haichen Shen, Meghan Cowan, Leyuan Wang, Yuwei Hu, Luis Ceze, et al. 2018. {TVM}: An automated {End-to-End} optimizing compiler for deep learning. In *13th USENIX Symposium on Operating Systems Design and Implementation (OSDI 18)*. 578–594.
- [64] Colin B. Clement, Dawn Drain, Jonathan Timcheck, Alexey Svyatkovskiy, and Neel Sundaresan. 2020. PyMT5: Multi-Mode Translation of Natural Language and Python Code with Transformers. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP*. Association for Computational Linguistics, 9052–9065.
- [65] Google Cloud. 2025. AutoML Solutions. <https://cloud.google.com/automl>. Accessed: 2025-09-15.
- [66] Google Cloud. 2025. Google Cloud CLI Documentation. <https://cloud.google.com/cli>. Accessed: 2025-09-15.
- [67] Tristan Coignon, Clément Quinton, and Romain Rouvoy. 2024. A Performance Study of LLM-Generated Code on Leetcode. In *International Conference on Evaluation and Assessment in Software Engineering (EASE)*. 79–89.
- [68] Antonia Creswell and Murray Shanahan. 2022. Faithful reasoning using large language models. *arXiv preprint arXiv:2208.14271* (2022).
- [69] Chris Cummins, Volker Seeker, Dejan Grubisic, Baptiste Roziere, Jonas Gehring, Gabriel Synnaeve, and Hugh Leather. 2025. LLM Compiler: Foundation Language Models for Compiler Optimization. In *Proceedings of the 34th ACM*

- SIGPLAN International Conference on Compiler Construction* (Las Vegas, NV, USA) (CC '25). Association for Computing Machinery, New York, NY, USA, 141–153. doi:10.1145/3708493.3712691
- [70] Inc. Cursor. 2023. Cursor: The AI Code Editor. Accessed: 2025-09-15.
 - [71] Mehul Damani, Idan Shenfeld, Andi Peng, Andreea Bobu, and Jacob Andreas. 2024. Learning how hard to think: Input-adaptive allocation of lm computation. *arXiv preprint arXiv:2410.04707* (2024).
 - [72] Edoardo Debenedetti, Jie Zhang, Mislav Balunović, Luca Beurer-Kellner, Marc Fischer, and Florian Tramèr. 2024. AgentDojo: A Dynamic Environment to Evaluate Prompt Injection Attacks and Defenses for LLM Agents. arXiv:2406.13352 [cs.CR] <https://arxiv.org/abs/2406.13352>
 - [73] Gelei Deng, Yi Liu, Victor Mayoral-Vilches, Peng Liu, Yuekang Li, Yuan Xu, Tianwei Zhang, Yang Liu, Martin Pinzger, and Stefan Rass. 2024. PentestGPT: An LLM-empowered Automatic Penetration Testing Tool. arXiv:2308.06782 [cs.SE] <https://arxiv.org/abs/2308.06782>
 - [74] Darshan Deshpande, Varun Gangal, Hersh Mehta, Jitin Krishnan, Anand Kannappan, and Rebecca Qian. 2025. TRAIL: Trace Reasoning and Agentic Issue Localization. arXiv:2505.08638 [cs.AI] <https://arxiv.org/abs/2505.08638>
 - [75] devgpt-labs. 2024. DevGPT: Transform Jira tickets straight into code. <https://github.com/devgpt-labs/devgpt>. GitHub repository. Accessed: 2025-08-31.
 - [76] Jacob Devlin, Jonathan Uesato, Surya Bhupatiraju, Rishabh Singh, Abdel-rahman Mohamed, and Pushmeet Kohli. 2017. Robustfill: Neural program learning under noisy i/o. In *International conference on machine learning*. PMLR, 990–998.
 - [77] Jiayu Ding, Shuming Ma, Li Dong, Xingxing Zhang, Shaohan Huang, Wenhui Wang, Nanning Zheng, and Furu Wei. 2023. Longnet: Scaling transformers to 1,000,000,000 tokens. *arXiv preprint arXiv:2307.02486* (2023).
 - [78] Yiran Ding, Li Lina Zhang, Chengruidong Zhang, Yuanyuan Xu, Ning Shang, Jiahang Xu, Fan Yang, and Mao Yang. 2024. Longrope: Extending llm context window beyond 2 million tokens. *arXiv preprint arXiv:2402.13753* (2024).
 - [79] Zhihua Duan and Jialin Wang. 2024. Exploration of llm multi-agent application implementation based on langgraph+crewai. *arXiv preprint arXiv:2411.18241* (2024).
 - [80] Shubhang Shekhar Dvivedi, Vyshnav Vijay, Sai Leela Rahul Pujari, Shoumik Lodh, and Dhruv Kumar. 2024. A comparative analysis of large language models for code documentation generation. In *Proceedings of the 1st ACM international conference on AI-powered software*. 65–73.
 - [81] Zachary Englhardt, Richard Li, Dilini Nissanka, Zhihan Zhang, Girish Narayanswamy, Joseph Breda, Xin Liu, Shwetak Patel, and Vikram Iyer. 2024. Exploring and characterizing large language models for embedded system development and debugging. In *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems*. 1–9.
 - [82] Gheorghe Comanici et al. 2025. Gemini 2.5: Pushing the Frontier with Advanced Reasoning, Multimodality, Long Context, and Next Generation Agentic Capabilities. arXiv:2507.06261 [cs.CL] <https://arxiv.org/abs/2507.06261>
 - [83] Sarah Fakhoury, Aaditya Naik, Georgios Sakkas, Saikat Chakraborty, and Shuvendu K Lahiri. 2024. Llm-based test-driven interactive code generation: User study and empirical evaluation. *IEEE Transactions on Software Engineering* (2024).
 - [84] Angela Fan, Beliz Gokkaya, Mark Harman, Mitya Lyubarskiy, Shubho Sengupta, Shin Yoo, and Jie M. Zhang. 2023. Large Language Models for Software Engineering: Survey and Open Problems. In *International Conference on Software Engineering: Future of Software Engineering, ICSE-FoSE 2023*. IEEE, 31–53.
 - [85] Zhiyu Fan, Xiang Gao, Martin Mirchev, Abhik Roychoudhury, and Shin Hwei Tan. 2023. Automated repair of programs from large language models. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*. IEEE, 1469–1481.
 - [86] Matthias Feurer, Aaron Klein, Katharina Eggenberger, Jost Springenberg, Manuel Blum, and Frank Hutter. 2015. Efficient and robust automated machine learning. *Advances in neural information processing systems* 28 (2015).
 - [87] Yao Fu, Leyang Xue, Yeqi Huang, Andrei-Octavian Brabete, Dmitrii Ustiugov, Yuvraj Patel, and Luo Mai. 2024. {ServerlessLLM}:{Low-Latency} serverless inference for large language models. In *18th USENIX Symposium on Operating Systems Design and Implementation (OSDI 24)*. 135–153.
 - [88] Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yixin Dai, Jiawei Sun, Haofen Wang, and Haofen Wang. 2023. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997* 2, 1 (2023).
 - [89] GitHub. 2021. Introducing GitHub Copilot: your AI pair programmer. <https://github.blog/2021-06-29-introducing-github-copilot-ai-pair-programmer/>
 - [90] GitHub. 2025. Understanding GitHub Actions. <https://docs.github.com/articles/getting-started-with-github-actions>. Accessed: 2025-09-15.
 - [91] Jingzhi Gong and Tao Chen. 2024. Deep Configuration Performance Learning: A Systematic Survey and Taxonomy. *ACM Transactions on Software Engineering and Methodology (TOSEM)* (2024).
 - [92] Jingzhi Gong, Vardan Voskanyan, Paul Brookes, Fan Wu, Wei Jie, Jie Xu, Rafail Giavrimis, Mike Basios, Leslie Kanthan, and Zheng Wang. 2025. Language models for code optimization: Survey, challenges and future directions. *arXiv*

- preprint *arXiv:2501.01277* (2025).
- [93] Google. 2015. gRPC: A High Performance, Open-Source Universal RPC Framework. <https://grpc.io/>.
 - [94] Sumit Gulwani, Oleksandr Polozov, Rishabh Singh, et al. 2017. Program synthesis. *Foundations and Trends® in Programming Languages* 4, 1-2 (2017), 1–119.
 - [95] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948* (2025).
 - [96] Francesco Guzzi, George Kourousias, Roberto Pugliese, Alessandra Gianoncelli, and Fulvio Billè. 2024. Translating and Optimising Computational Microscopy Algorithms with Large Language Models. In *47th MIPRO ICT and Electronics Convention, MIPRO*. IEEE, 1700–1705.
 - [97] Jo E Hannay, Tore Dybå, Erik Arisholm, and Dag IK Sjøberg. 2009. The effectiveness of pair programming: A meta-analysis. *Information and software technology* 51, 7 (2009), 1110–1122.
 - [98] Jingxuan He, Mark Vero, Gabriela Krasnopolska, and Martin Vechev. 2024. Instruction tuning for secure code generation. *arXiv preprint arXiv:2402.09497* (2024).
 - [99] Xin He, Kaiyong Zhao, and Xiaowen Chu. 2021. AutoML: A survey of the state-of-the-art. *Knowledge-based systems* 212 (2021), 106622.
 - [100] Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, et al. 2023. Metagpt: Meta programming for multi-agent collaborative framework. *arXiv preprint arXiv:2308.00352* 3, 4 (2023), 6.
 - [101] Coleman Hooper, Sehoon Kim, Hiva Mohammadzadeh, Michael W Mahoney, Yakun S Shao, Kurt Keutzer, and Amir Gholami. 2024. Kvquant: Towards 10 million context length llm inference with kv cache quantization. *Advances in Neural Information Processing Systems* 37 (2024), 1270–1303.
 - [102] Peyman Hosseini, Ignacio Castro, Iacopo Ghinassi, and Matthew Purver. 2024. Efficient solutions for an intriguing failure of llms: Long context window does not mean llms can analyze long sequences flawlessly. *arXiv preprint arXiv:2408.01866* (2024).
 - [103] Xinyi Hou, Yanjie Zhao, Yue Liu, Zhou Yang, Kailong Wang, Li Li, Xiapu Luo, David Lo, John Grundy, and Haoyu Wang. 2024. Large language models for software engineering: A systematic literature review. *ACM Transactions on Software Engineering and Methodology* 33, 8 (2024), 1–79.
 - [104] Wenyue Hua, Xianjun Yang, Mingyu Jin, Zelong Li, Wei Cheng, Ruixiang Tang, and Yongfeng Zhang. 2024. Trustagent: Towards safe and trustworthy llm-based agents. *arXiv preprint arXiv:2402.01586* (2024).
 - [105] Dong Huang, Jianbo Dai, Han Weng, Puzhen Wu, QING Yuhao, Heming Cui, Zhijiang Guo, and Jie Zhang. 2024. EffiLearner: Enhancing Efficiency of Generated Code via Self-Optimization. (2024).
 - [106] Dong Huang, Jie M Zhang, Qingwen Bu, Xiaofei Xie, Junjie Chen, and Heming Cui. 2024. Bias testing and mitigation in llm-based code generation. *ACM Transactions on Software Engineering and Methodology* (2024).
 - [107] Dong Huang, Jie M Zhang, Michael Luck, Qingwen Bu, Yuhao Qing, and Heming Cui. 2023. Agentcoder: Multi-agent-based code generation with iterative testing and optimisation. *arXiv preprint arXiv:2312.13010* (2023).
 - [108] Jie Huang and Kevin Chen-Chuan Chang. 2023. Citation: A key to building responsible and accountable large language models. *arXiv preprint arXiv:2307.02185* (2023).
 - [109] Xu Huang, Weiwen Liu, Xiaolong Chen, Xingmei Wang, Hao Wang, Defu Lian, Yasheng Wang, Ruiming Tang, and Enhong Chen. 2024. Understanding the planning of LLM agents: A survey. *arXiv preprint arXiv:2402.02716* (2024).
 - [110] Hugging Face. 2025. Hugging Face. <https://huggingface.co/>. Accessed: 2025-09-11.
 - [111] Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, et al. 2024. Qwen2.5-coder Technical Report. *arXiv:2409.12186* (2024).
 - [112] Jez Humble and David Farley. 2010. *Continuous delivery: reliable software releases through build, test, and deployment automation*. Pearson Education.
 - [113] Rasha Ahmad Husein, Hala Aburajouh, and Catagay Catal. 2025. Large language models for code completion: A systematic literature review. *Computer Standards & Interfaces* 92 (2025), 103917.
 - [114] GitLab Inc. 2025. Get started with GitLab CI/CD. <https://docs.gitlab.com/ci/>. Accessed: 2025-09-15.
 - [115] Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. 2024. Livecodebench: Holistic and contamination free evaluation of large language models for code. *arXiv preprint arXiv:2403.07974* (2024).
 - [116] Naman Jain, Skanda Vaidyanath, Arun Iyer, Nagarajan Natarajan, Suresh Parthasarathy, Sriram Rajamani, and Rahul Sharma. 2022. Jigsaw: Large Language Models Meet Program Synthesis. In *International Conference on Software Engineering (ICSE)*. 1219–1231.
 - [117] Nicholas R Jennings. 2000. On agent-based software engineering. *Artificial intelligence* 117, 2 (2000), 277–296.
 - [118] Huiqiang Jiang, Qianhui Wu, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. 2023. LlmLingua: Compressing prompts for accelerated inference of large language models. *arXiv preprint arXiv:2310.05736* (2023).

- [119] Juyong Jiang, Fan Wang, Jiasi Shen, Sungju Kim, and Sunghun Kim. 2025. A Survey on Large Language Models for Code Generation. *ACM Trans. Softw. Eng. Methodol.* (July 2025). doi:10.1145/3747588 Just Accepted.
- [120] Xue Jiang, Yihong Dong, Yongding Tao, Huanyu Liu, Zhi Jin, Wenpin Jiao, and Ge Li. 2024. Rocode: Integrating backtracking mechanism and program analysis in large language models for code generation. *arXiv preprint arXiv:2411.07112* (2024).
- [121] Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. 2023. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770* (2023).
- [122] Sathvik Joel, Jie JW Wu, and Fatemeh H Fard. 2024. Survey on Code Generation for Low resource and Domain Specific Programming Languages. *arXiv:2410.03981* (2024).
- [123] Steven Jorgensen, Giorgia Nadizar, Gloria Pietropoli, Luca Manzoni, Eric Medvet, Una-May O'Reilly, and Erik Hemberg. 2024. Large language model-based test case generation for gp agents. In *Proceedings of the genetic and evolutionary computation conference*. 914–923.
- [124] GIREESH KAMBALA. 2024. Intelligent Software Agents for Continuous Delivery: Leveraging AI and Machine Learning for Fully Automated DevOps Pipelines. *Iconic Research And Engineering Journals* 8, 1 (2024), 662–670.
- [125] Sungmin Kang, Bei Chen, Shin Yoo, and Jian-Guang Lou. 2025. Explainable automated debugging via large language model-driven scientific debugging. *Empirical Software Engineering* 30, 2 (2025), 45.
- [126] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361* (2020).
- [127] Sehoon Kim, Suhong Moon, Ryan Tabrizi, Nicholas Lee, Michael W Mahoney, Kurt Keutzer, and Amir Gholami. 2024. An llm compiler for parallel function calling. In *Forty-first International Conference on Machine Learning*.
- [128] Barbara Kitchenham and Stuart M. Charters. 2007. Guidelines for Performing Systematic Literature Reviews in Software Engineering.
- [129] KServe Community. 2025. KServe: Serverless Inferencing on Kubernetes. <https://kserve.github.io/website/>. Accessed: 2025-09-11.
- [130] Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2023. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation. *arXiv preprint arXiv:2302.09664* (2023).
- [131] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the 29th symposium on operating systems principles*. 611–626.
- [132] GitHub Security Lab. 2020. CodeQL: Semantic Code Analysis Engine. <https://codeql.github.com/>.
- [133] Chris Lattner and Vikram Adve. 2004. LLVM: A Compilation Framework for Lifelong Program Analysis & Transformation. In *Proceedings of the International Symposium on Code Generation and Optimization: Feedback-Directed and Runtime Optimization* (Palo Alto, California) (CGO '04). IEEE Computer Society, USA, 75.
- [134] Chris Lattner and Vikram Adve. 2004. LLVM: A compilation framework for lifelong program analysis & transformation. In *International symposium on code generation and optimization, 2004. CGO 2004*. IEEE, 75–86.
- [135] Chris Lattner, Mehdi Amini, Uday Bondhugula, Albert Cohen, Andy Davis, Jacques Pienaar, River Riddle, Tatiana Shpeisman, Nicolas Vasilache, and Oleksandr Zinenko. 2021. MLIR: Scaling compiler infrastructure for domain specific computation. In *2021 IEEE/ACM International Symposium on Code Generation and Optimization (CGO)*. IEEE, 2–14.
- [136] Hung Le, Yue Wang, Akhilesh Deepak Gotmare, Silvio Savarese, and Steven Chu-Hong Hoi. 2022. CodeRL: Mastering Code Generation through Pretrained Models and Deep Reinforcement Learning. In *Neural Information Processing Systems (NeurIPS)*.
- [137] Thanh Le-Cong, Bach Le, and Toby Murray. 2025. Can LLMs Reason About Program Semantics? A Comprehensive Evaluation of LLMs on Formal Specification Inference. *arXiv:2503.04779* [cs.PL] <https://arxiv.org/abs/2503.04779>
- [138] Wonbeom Lee, Jungi Lee, Junghwan Seo, and Jaewoong Sim. 2024. {InfiniGen}: Efficient generative inference of large language models with dynamic {KV} cache management. In *18th USENIX Symposium on Operating Systems Design and Implementation (OSDI 24)*. 155–172.
- [139] Fangyu Lei, Jixuan Chen, Yuxiao Ye, Ruisheng Cao, Dongchan Shin, Hongjin Su, Zhaoqing Suo, Hongcheng Gao, Wenjing Hu, Pengcheng Yin, et al. 2024. Spider 2.0: Evaluating language models on real-world enterprise text-to-sql workflows. *arXiv preprint arXiv:2411.07763* (2024).
- [140] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in neural information processing systems* 33 (2020), 9459–9474.
- [141] Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2023. Camel: Communicative agents for" mind" exploration of large language model society. *Advances in Neural Information Processing Systems* 36 (2023), 51991–52008.

- [142] Jia Li, Ge Li, Yongmin Li, and Zhi Jin. 2025. Structured chain-of-thought prompting for code generation. *ACM Transactions on Software Engineering and Methodology* 34, 2 (2025), 1–23.
- [143] Xinzhe Li. 2025. A review of prominent paradigms for llm-based agents: Tool use, planning (including rag), and feedback learning. In *Proceedings of the 31st International Conference on Computational Linguistics*. 9760–9779.
- [144] Xiaoxi Li, Jiajie Jin, Yujia Zhou, Yuyao Zhang, Peitian Zhang, Yutao Zhu, and Zhicheng Dou. 2025. From matching to generation: A survey on generative information retrieval. *ACM Transactions on Information Systems* 43, 3 (2025), 1–62.
- [145] Xinyi Li, Sai Wang, Siqi Zeng, Yu Wu, and Yi Yang. 2024. A survey on LLM-based multi-agent systems: workflow, infrastructure, and challenges. *Vicinagearth* 1, 1 (2024), 9.
- [146] Yujia Li, David Choi, Junyoung Chung, and et al. 2022. Competition-Level Code Generation with AlphaCode. *Science* 378 (2022), 1092–1097.
- [147] Ziyang Li, Saikat Dutta, and Mayur Naik. 2024. IRIS: LLM-assisted static analysis for detecting security vulnerabilities. *arXiv preprint arXiv:2405.17238* (2024).
- [148] Bo Lin, Shangwen Wang, Yihao Qin, Liqian Chen, and Xiaoguang Mao. 2025. Give LLMs a Security Course: Securing Retrieval-Augmented Code Generation via Knowledge Injection. *arXiv:2504.16429* [cs.CR] <https://arxiv.org/abs/2504.16429>
- [149] Fang Liu, Yang Liu, Lin Shi, Houkun Huang, Ruifeng Wang, Zhen Yang, Li Zhang, Zhongqi Li, and Yuchi Ma. 2024. Exploring and evaluating hallucinations in llm-powered code generation. *arXiv preprint arXiv:2404.00971* (2024).
- [150] Junwei Liu, Kaixin Wang, Yixuan Chen, Xin Peng, Zhenpeng Chen, Lingming Zhang, and Yiling Lou. 2024. Large Language Model-Based Agents for Software Engineering: A Survey. *arXiv:2409.02977* (2024).
- [151] Kaiyuan Liu, Youcheng Pan, Yang Xiang, Daojing He, Jing Li, Yexing Du, and Tianrun Gao. 2025. ProjectEval: A Benchmark for Programming Agents Automated Evaluation on Project-Level Code Generation. In *Findings of the Association for Computational Linguistics: ACL 2025*, Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (Eds.). Association for Computational Linguistics, Vienna, Austria, 20205–20221. doi:10.18653/v1/2025.findings-acl.1036
- [152] Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023. Lost in the middle: How language models use long contexts. *arXiv preprint arXiv:2307.03172* (2023).
- [153] Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, Shudan Zhang, Xiang Deng, Aohan Zeng, Zhengxiao Du, Chenhui Zhang, Sheng Shen, Tianjun Zhang, Yu Su, Huan Sun, Minlie Huang, Yuxiao Dong, and Jie Tang. 2023. AgentBench: Evaluating LLMs as Agents. *arXiv preprint arXiv: 2308.03688* (2023).
- [154] Yan Liu, Xiaokang Chen, Yan Gao, Zhe Su, Fengji Zhang, Daoguang Zan, Jian-Guang Lou, Pin-Yu Chen, and Tsung-Yi Ho. 2023. Uncovering and quantifying social biases in code generation. *Advances in Neural Information Processing Systems* 36 (2023), 2368–2380.
- [155] Yuhao Liu, Hanchen Li, Yihua Cheng, Siddhant Ray, Yuyang Huang, Qizheng Zhang, Kuntai Du, Jiayi Yao, Shan Lu, Ganesh Ananthanarayanan, et al. 2024. CacheGen: Kv cache compression and streaming for fast large language model serving. In *Proceedings of the ACM SIGCOMM 2024 Conference*. 38–56.
- [156] Yue Liu, Chakkrit Tantithamthavorn, Yonghui Liu, and Li Li. 2024. On the reliability and explainability of language models for program generation. *ACM Transactions on Software Engineering and Methodology* 33, 5 (2024), 1–26.
- [157] Zhiwei Liu, Weiran Yao, Jianguo Zhang, Le Xue, Shelby Heinecke, Rithesh Murthy, Yihao Feng, Zeyuan Chen, Juan Carlos Nieves, Devansh Arpit, et al. 2023. Bolaa: Benchmarking and orchestrating llm-augmented autonomous agents. *arXiv preprint arXiv:2308.05960* (2023).
- [158] Zirui Liu, Jiayi Yuan, Hongye Jin, Shaochen Zhong, Zhaozhuo Xu, Vladimir Braverman, Beidi Chen, and Xia Hu. 2024. Kivi: A tuning-free asymmetric 2bit quantization for kv cache. *arXiv preprint arXiv:2402.02750* (2024).
- [159] Chloe Loughridge, Qinyi Sun, Seth Ahrenbach, Federico Cassano, Chuyue Sun, Ying Sheng, Anish Mudide, Md Rakib Hossain Misu, Nada Amin, and Max Tegmark. 2024. DafnyBench: A Benchmark for Formal Software Verification. *arXiv preprint arXiv:2406.08467* (2024).
- [160] Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, et al. 2024. StarCoder 2 and The Stack v2: The Next Generation. *arXiv:2402.19173* (2024).
- [161] Aman Madaan, Niket Tandon, Prakhar Gupta, and et al. 2023. Self-Refine: Iterative Refinement with Self-Feedback. In *Neural Information Processing Systems (NeurIPS)*.
- [162] Adyasha Maharana, Dong-Ho Lee, Sergey Tulyakov, Mohit Bansal, Francesco Barbieri, and Yuwei Fang. 2024. Evaluating Very Long-Term Conversational Memory of LLM Agents. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 13851–13870.
- [163] Saeed Maleki, Yaoqing Gao, Maria J Garzar, Tommy Wong, David A Padua, et al. 2011. An Evaluation of Vectorizing Compilers. In *2011 International Conference on Parallel Architectures and Compilation Techniques*. IEEE, 372–382.

- [164] Saeed Maleki, Yaoqing Gao, María J. Garzarán, Tommy Wong, and David A. Padua. 2011. An Evaluation of Vectorizing Compilers. In *2011 International Conference on Parallel Architectures and Compilation Techniques*. 372–382. doi:10.1109/PACT.2011.68
- [165] Valentin JM Manès, HyungSeok Han, Choongwoo Han, Sang Kil Cha, Manuel Egele, Edward J Schwartz, and Maverick Woo. 2019. The art, science, and engineering of fuzzing: A survey. *IEEE Transactions on Software Engineering* 47, 11 (2019), 2312–2331.
- [166] Rohin Manvi, Anikait Singh, and Stefano Ermon. 2024. Adaptive inference-time compute: Llms can predict if they can do better, even mid-generation. *arXiv preprint arXiv:2410.02725* (2024).
- [167] Mark Marron. 2023. Programming Languages for AI Programing Agents (Invited Talk). In *Proceedings of the 19th ACM SIGPLAN International Symposium on Dynamic Languages (Cascais, Portugal) (DLS 2023)*. Association for Computing Machinery, New York, NY, USA, 7. doi:10.1145/3622759.3628225
- [168] Meta. 2024. Introducing Llama 3.1: Our most capable models to date. <https://ai.meta.com/blog/meta-llama-3-1/>
- [169] Microsoft. 2025. Azure Machine Learning. <https://azure.microsoft.com/en-us/products/machine-learning/>. Accessed: 2025-09-11.
- [170] Philipp Moritz, Robert Nishihara, Stephanie Wang, Alexey Tumanov, Richard Liaw, Eric Liang, Melih Elibol, Zongheng Yang, William Paul, Michael I Jordan, et al. 2018. Ray: A distributed framework for emerging {AI} applications. In *13th USENIX symposium on operating systems design and implementation (OSDI 18)*. 561–577.
- [171] Fangwen Mu, Lin Shi, Song Wang, Zhuohao Yu, Binquan Zhang, ChenXue Wang, Shichao Liu, and Qing Wang. 2024. Clarifygpt: A framework for enhancing llm-based code generation via requirements clarification. *Proceedings of the ACM on Software Engineering* 1, FSE (2024), 2332–2354.
- [172] Nitin Naik and Paul Jenkins. 2016. Web protocols and challenges of web latency in the web of things. In *2016 Eighth International Conference on ubiquitous and future networks (ICUFN)*. IEEE, 845–850.
- [173] Oscar Nierstrasz. 1993. Regular types for active objects. *ACM sigplan Notices* 28, 10 (1993), 1–15.
- [174] Oscar M Nierstrasz. 1987. Active objects in Hybrid. In *Conference proceedings on Object-oriented programming systems, languages and applications*. 243–253.
- [175] Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan Wang, Yingbo Zhou, Silvio Savarese, and Caiming Xiong. 2023. CodeGen: An Open Large Language Model for Code with Multi-Turn Program Synthesis. In *The Eleventh International Conference on Learning Representations, ICLR*. OpenReview.net.
- [176] Alexander Novikov, Ngàn Vũ, Marvin Eisenberger, Emilien Dupont, Po-Sen Huang, Adam Zsolt Wagner, Sergey Shirobokov, Borislav Kozlovskii, Francisco JR Ruiz, Abbas Mehrabian, et al. 2025. AlphaEvolve: A coding agent for scientific and algorithmic discovery. *arXiv preprint arXiv:2506.13131* (2025).
- [177] Ana Nunez, Nafis Tanveer Islam, Sumit Kumar Jha, and Peyman Najafirad. 2024. Autosafecoder: A multi-agent framework for securing llm code generation through static analysis and fuzz testing. *arXiv preprint arXiv:2409.10737* (2024).
- [178] NVIDIA Technical Blog. 2025. Safeguard Agentic AI Systems with the NVIDIA Safety Recipe. <https://developer.nvidia.com/blog/safeguard-agentic-ai-systems-with-the-nvidia-safety-recipe>.
- [179] OpenAI. 2021. Evaluating Large Language Models Trained on Code. (2021). arXiv:2107.03374
- [180] OpenAI. 2024. Hello GPT-4o. <https://openai.com/index/hello-gpt-4o/>
- [181] OpenAI. 2025. Introducing GPT-5. <https://openai.com/index/introducing-gpt-5/>. Accessed: 2025-08-10.
- [182] OpenFOAM Foundation. 2023. *OpenFOAM: The Open Source CFD Toolbox*. <https://www.openfoam.com>.
- [183] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems* 35 (2022), 27730–27744.
- [184] Charles Packer, Vivian Fang, Shishir_G Patil, Kevin Lin, Sarah Wooders, and Joseph_E Gonzalez. 2023. MemGPT: Towards LLMs as Operating Systems. (2023).
- [185] Pratyush Patel, Esha Choukse, Chaojie Zhang, Aashaka Shah, Íñigo Goiri, Saeed Maleki, and Ricardo Bianchini. 2024. Splitwise: Efficient generative llm inference using phase splitting. In *2024 ACM/IEEE 51st Annual International Symposium on Computer Architecture (ISCA)*. IEEE, 118–132.
- [186] Shishir G Patil, Tianjun Zhang, Vivian Fang, Roy Huang, Aaron Hao, Martin Casado, Joseph E Gonzalez, Raluca Ada Popa, Ion Stoica, et al. 2024. Goex: Perspectives and designs towards a runtime for autonomous llm applications. *arXiv preprint arXiv:2404.06921* (2024).
- [187] Hammond Pearce, Baleegh Ahmad, Benjamin Tan, Brendan Dolan-Gavitt, and Ramesh Karri. 2025. Asleep at the keyboard? assessing the security of github copilot’s code contributions. *Commun. ACM* 68, 2 (2025), 96–105.
- [188] Yun Peng, Akhilesh Deepak Gotmare, Michael Lyu, Caiming Xiong, Silvio Savarese, and Doyen Sahoo. 2024. Perf-CodeGen: Improving Performance of LLM Generated Code with Execution Feedback. *arXiv:2412.03578* (2024).
- [189] Cody J Permann, Derek R Gaston, David Andrš, Robert W Carlsen, Fande Kong, Alexander D Lindsay, Jason M Miller, John W Peterson, Andrew E Slaughter, Roy H Stogner, et al. 2020. MOOSE: Enabling massively parallel multiphysics

- simulation. *SoftwareX* 11 (2020), 100430.
- [190] The Jenkins Project. 2025. Jenkins (software). [https://en.wikipedia.org/wiki/Jenkins_\(software\)](https://en.wikipedia.org/wiki/Jenkins_(software)). Accessed: 2025-09-15.
 - [191] Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan Dang, Jiahao Li, Cheng Yang, Weize Chen, Yusheng Su, Xin Cong, Juyuan Xu, Dahai Li, Zhiyuan Liu, and Maosong Sun. 2024. ChatDev: Communicative Agents for Software Development. arXiv:2307.07924 [cs.SE] <https://arxiv.org/abs/2307.07924>
 - [192] Hongjin Qian, Zheng Liu, Peitian Zhang, Kelong Mao, Defu Lian, Zhicheng Dou, and Tiejun Huang. 2025. Memorag: Boosting long context processing with global memory-enhanced retrieval augmentation. In *Proceedings of the ACM on Web Conference 2025*. 2366–2377.
 - [193] Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Lauren Hong, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein, Dahai Li, Zhiyuan Liu, and Maosong Sun. 2023. ToolLLM: Facilitating Large Language Models to Master 16000+ Real-world APIs. arXiv:2307.16789 [cs.AI] <https://arxiv.org/abs/2307.16789>
 - [194] Yuhao Qing, Boyu Zhu, Mingzhe Du, Zhijiang Guo, Terry Yue Zhuo, Qianru Zhang, Jie M. Zhang, Heming Cui, Siu-Ming Yiu, Dong Huang, See-Kiong Ng, and Luu Anh Tuan. 2025. EffiBench-X: A Multi-Language Benchmark for Measuring Efficiency of LLM-Generated Code. arXiv:2505.13004 [cs.CL] <https://arxiv.org/abs/2505.13004>
 - [195] Umair Qudus, Michael Röder, Muhammad Saleem, and Axel-Cyrille Ngonga Ngomo. 2025. Fact Checking Knowledge Graphs—A Survey. *Comput. Surveys* (2025).
 - [196] Nazneen Fatema Rajani, Bryan McCann, Caiming Xiong, and Richard Socher. 2019. Explain yourself! leveraging language models for commonsense reasoning. *arXiv preprint arXiv:1906.02361* (2019).
 - [197] Arthur J Riel. 1996. *Object-oriented design heuristics*. Addison-Wesley Longman Publishing Co., Inc.
 - [198] Joel Rorseth, Parke Godfrey, Lukasz Golab, Divesh Srivastava, and Jarek Szlichta. 2025. LADYBUG: an LLM Agent DeBUGger for data-driven applications. (2025).
 - [199] Ranjan Sapkota, Konstantinos I. Roumeliotis, and Manoj Karkee. 2025. AI Agents vs. Agentic AI: A Conceptual Taxonomy, Applications and Challenges. arXiv:2505.10468 [cs.AI] <https://arxiv.org/abs/2505.10468>
 - [200] Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems* 36 (2023), 68539–68551.
 - [201] Sander Schulhoff, Michael Ilie, Nishant Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si, Yinheng Li, Aayush Gupta, Hyojung Han, Sevien Schulhoff, et al. 2024. The prompt report: a systematic survey of prompt engineering techniques. *arXiv preprint arXiv:2406.06608* (2024).
 - [202] Amazon Web Services. 2025. Amazon SageMaker Autopilot. <https://aws.amazon.com/sagemaker-ai/autopilot/>. Accessed: 2025-09-15.
 - [203] Chirag Shah. 2025. From Prompt Engineering to Prompt Science with Humans in the Loop. *Commun. ACM* 68, 6 (2025), 54–61.
 - [204] Manish Shetty, Naman Jain, Jinjian Liu, Vijay Kethanaboyina, Koushik Sen, and Ion Stoica. 2025. GSO: Challenging Software Optimization Tasks for Evaluating SWE-Agents. *arXiv preprint arXiv:2505.23671* (2025).
 - [205] Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: Language Agents with Verbal Reinforcement Learning. In *Neural Information Processing Systems (NeurIPS)*.
 - [206] Sergi Siso, Wes Armour, and Jeyarajan Thiayagalingam. 2019. Evaluating auto-vectorizing compilers through objective withdrawal of useful information. *ACM Transactions on Architecture and Code Optimization (TACO)* 16, 4 (2019), 1–23.
 - [207] Armando Solar-Lezama. 2008. *Program synthesis by sketching*. University of California, Berkeley.
 - [208] Zirui Song, Bin Yan, Yuhan Liu, Miao Fang, Mingzhe Li, Rui Yan, and Xiuying Chen. 2025. Injecting domain-specific knowledge into large language models: a comprehensive survey. *arXiv preprint arXiv:2502.10708* (2025).
 - [209] Weihang Su, Yichen Tang, Qingyao Ai, Zhijing Wu, and Yiqun Liu. 2024. DRAGIN: dynamic retrieval augmented generation based on the information needs of large language models. *arXiv preprint arXiv:2403.10081* (2024).
 - [210] Chuanneng Sun, Songjun Huang, and Dario Pompili. 2024. Llm-based multi-agent reinforcement learning: Current and future directions. *arXiv preprint arXiv:2405.11106* (2024).
 - [211] Hanshi Sun, Li-Wen Chang, Wenlei Bao, Size Zheng, Ningxin Zheng, Xin Liu, Harry Dong, Yuejie Chi, and Beidi Chen. 2024. Shadowkv: Kv cache in shadows for high-throughput long-context llm inference. *arXiv preprint arXiv:2410.21465* (2024).
 - [212] Haotian Sun, Yuchen Zhuang, Ling kai Kong, Bo Dai, and Chao Zhang. 2023. Adapllanner: Adaptive planning from feedback with language models. *Advances in neural information processing systems* 36 (2023), 58202–58245.
 - [213] Jiao Sun, Q Vera Liao, Michael Muller, Mayank Agarwal, Stephanie Houde, Kartik Talamadupula, and Justin D Weisz. 2022. Investigating explainability of generative AI for code through scenario-based design. In *Proceedings of the 27th international conference on intelligent user interfaces*. 212–228.

- [214] Tianxiang Sun, Zhengfu He, Qin Zhu, Xipeng Qiu, and Xuanjing Huang. 2023. Multitask Pre-training of Modular Prompt for Chinese Few-Shot Learning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (Eds.). Association for Computational Linguistics, Toronto, Canada, 11156–11172. doi:10.18653/v1/2023.acl-long.625
- [215] Tabnine. 2025. Tabnine: AI coding assistant. <https://www.tabnine.com/>
- [216] Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530* (2024).
- [217] Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen, Yanru Chen, Yuankun Chen, Yutian Chen, et al. 2025. Kimi K2: Open Agentic Intelligence. *arXiv preprint arXiv:2507.20534* (2025).
- [218] The Terminal-Bench Team. 2025. Terminal-Bench: A Benchmark for AI Agents in Terminal Environments. <https://github.com/laude-institute/terminal-bench>
- [219] The MathWorks, Inc. 2023. *Simulink: Dynamic System Simulation for MATLAB*. <https://www.mathworks.com/products/simulink.html>.
- [220] Richard Jiarui Tong, Haoyang Li, Sridhar Raghavan, Qingsong Wen, Shannon Gray, Anand Paul, Joleen Liang, Janusz Zalewski, Yacheng Yang, George Tambouratzis, et al. 2025. IEEE AI Standards for Agentic Systems. In *2025 IEEE Conference on Artificial Intelligence (CAI)*. IEEE, 1603–1609.
- [221] Harold Triedman, Rishi Jha, and Vitaly Shmatikov. 2025. Multi-Agent Systems Execute Arbitrary Malicious Code. arXiv:2503.12188 [cs.CR] <https://arxiv.org/abs/2503.12188>
- [222] Lukas Twist, Jie M Zhang, Mark Harman, Don Syme, Joost Noppen, and Detlef Nauck. 2025. LLMs Love Python: A Study of LLMs’ Bias for Programming Languages and Libraries. *arXiv preprint arXiv:2503.17181* (2025).
- [223] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems*. 5998–6008.
- [224] Abhishek Verma, Luis Pedrosa, Madhukar Korupolu, David Oppenheimer, Eric Tune, and John Wilkes. 2015. Large-scale cluster management at Google with Borg. In *Proceedings of the Tenth European Conference on Computer Systems (Bordeaux, France) (EuroSys ’15)*. Association for Computing Machinery, New York, NY, USA, Article 18, 17 pages. doi:10.1145/2741948.2741964
- [225] Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. 2023. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv:2305.16291* (2023).
- [226] Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, et al. 2024. A survey on large language model based autonomous agents. *Frontiers of Computer Science* 18, 6 (2024), 186345.
- [227] Simin Wang, Liguao Huang, Amiao Gao, Jidong Ge, Tengfei Zhang, Haitao Feng, Ishna Satyarth, Ming Li, He Zhang, and Vincent Ng. 2023. Machine/Deep Learning for Software Engineering: A Systematic Literature Review. *IEEE Trans. Software Eng.* 49, 3 (2023), 1188–1231.
- [228] Weizhi Wang, Li Dong, Hao Cheng, Xiaodong Liu, Xifeng Yan, Jianfeng Gao, and Furu Wei. 2023. Augmenting language models with long-term memory. *Advances in Neural Information Processing Systems* 36 (2023), 74530–74543.
- [229] Xingyao Wang, Yangyi Chen, Lifan Yuan, Yizhe Zhang, Yunzhu Li, Hao Peng, and Heng Ji. 2024. Executable Code Actions Elicit Better LLM Agents. arXiv:2402.01030 [cs.CL] <https://arxiv.org/abs/2402.01030>
- [230] Xingyao Wang, Boxuan Li, Yufan Song, Frank F Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan, Yueqi Song, Bowen Li, Jaskirat Singh, et al. 2024. Openhands: An open platform for ai software developers as generalist agents. *arXiv preprint arXiv:2407.16741* (2024).
- [231] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems* 35 (2022), 24824–24837.
- [232] Jiaxin Wen, Jian Guan, Hongning Wang, Wei Wu, and Minlie Huang. 2024. Codeplan: Unlocking reasoning potential in large language models by scaling code-form planning. In *The Thirteenth International Conference on Learning Representations*.
- [233] Richard Wollheim. 1984. *The thread of life*. Yale University Press.
- [234] Jim Woodcock, Peter Gorm Larsen, Juan Bicarregui, and John Fitzgerald. 2009. Formal methods: Practice and experience. *ACM computing surveys (CSUR)* 41, 4 (2009), 1–36.
- [235] Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, et al. 2024. Autogen: Enabling next-gen LLM applications via multi-agent conversations. In *First Conference on Language Modeling*.
- [236] xAI. 2025. Grok 4: The Next-Generation AI Model. <https://x.ai/news/grok-4/>. Accessed: 2025-09-15.

- [237] Frank F. Xu, Uri Alon, Graham Neubig, and Vincent Josua Hellendoorn. 2022. A Systematic Evaluation of Large Language Models of Code. In *SIGPLAN International Symposium on Machine Programming*. ACM, 1–10.
- [238] Jing Xu. 2008. Rule-based automatic software performance diagnosis and improvement. In *Proceedings of the 7th international workshop on Software and performance*. 1–12.
- [239] Jie Xu, Dacheng Zhang, Lu Liu, and Xianxian Li. 2012. Dynamic Authentication for Cross-Realm SOA-Based Business Processes. *IEEE Transactions on Services Computing* 5, 1 (2012), 20–32. doi:10.1109/TSC.2010.33
- [240] Kai Xu, YiWei Mao, XinYi Guan, and ZiLong Feng. 2025. Web-Bench: A LLM Code Benchmark Based on Web Standards and Frameworks. *arXiv preprint arXiv:2505.07473* (2025).
- [241] Dongjie Yang, XiaoDong Han, Yan Gao, Yao Hu, Shilin Zhang, and Hai Zhao. 2024. Pyramidinfer: Pyramid kv cache compression for high-throughput llm inference. *arXiv preprint arXiv:2405.12532* (2024).
- [242] John Yang, Carlos E Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan, and Ofir Press. 2024. Swe-agent: Agent-computer interfaces enable automated software engineering. *Advances in Neural Information Processing Systems* 37 (2024), 50528–50652.
- [243] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023. ReAct: Synergizing Reasoning and Acting in Language Models. In *International Conference on Learning Representations (ICLR)*. <https://arxiv.org/abs/2210.03629>
- [244] Gyeong-In Yu, Joo Seong Jeong, Geon-Woo Kim, Soojeong Kim, and Byung-Gon Chun. 2022. Orca: A distributed serving system for {Transformer-Based} generative models. In *16th USENIX Symposium on Operating Systems Design and Implementation (OSDI 22)*. 521–538.
- [245] Minchen Yu, Zhifeng Jiang, Hok Chun Ng, Wei Wang, Ruichuan Chen, and Bo Li. 2021. Gillis: Serving large neural networks in serverless functions with automatic model partitioning. In *2021 IEEE 41st International Conference on Distributed Computing Systems (ICDCS)*. IEEE, 138–148.
- [246] Eric Zelikman, Eliana Lorch, Lester Mackey, and Adam Tauman Kalai. 2024. Self-Taught Optimizer (STOP): Recursively Self-Improving Code Generation. (2024).
- [247] Qiusi Zhan, Zhixiang Liang, Zifan Ying, and Daniel Kang. 2024. InjecAgent: Benchmarking Indirect Prompt Injections in Tool-Integrated Large Language Model Agents. arXiv:2403.02691 [cs.CL] <https://arxiv.org/abs/2403.02691>
- [248] Michael JQ Zhang, W Bradley Knox, and Eunsol Choi. 2024. Modeling future conversation turns to teach llms to ask clarifying questions. *arXiv preprint arXiv:2410.13788* (2024).
- [249] Quanjun Zhang, Chunrong Fang, Yang Xie, Yuxiang Ma, Weisong Sun, Yun Yang, and Zhenyu Chen. 2024. A Systematic Literature Review on Large Language Models for Automated Program Repair. (2024). arXiv:2405.01466
- [250] Yuntong Zhang, Haifeng Ruan, Zhiyu Fan, and Abhik Roychoudhury. 2024. AutoCodeRover: Autonomous Program Improvement. In *ACM SIGSOFT International Symposium on Software Testing and Analysis*. 1592–1604.
- [251] Zhexin Zhang, Shiyao Cui, Yida Lu, Jingzhuo Zhou, Junxiao Yang, Hongning Wang, and Minlie Huang. 2025. Agent-SafetyBench: Evaluating the Safety of LLM Agents. arXiv:2412.14470 [cs.CL] <https://arxiv.org/abs/2412.14470>
- [252] Yilong Zhao, Chien-Yu Lin, Kan Zhu, Zihao Ye, Lequn Chen, Size Zheng, Luis Ceze, Arvind Krishnamurthy, Tianqi Chen, and Baris Kasikci. 2024. Atom: Low-bit quantization for efficient and accurate llm serving. *Proceedings of Machine Learning and Systems* 6 (2024), 196–209.
- [253] Dewu Zheng, Yanlin Wang, Ensheng Shi, Hongyu Zhang, and Zibin Zheng. 2024. How well do llms generate code for different application domains? benchmark and evaluation. *arXiv preprint arXiv:2412.18573* (2024).
- [254] Qinkai Zheng, Xiao Xia, Xu Zou, Yuxiao Dong, Shan Wang, Yufei Xue, Lei Shen, Zihan Wang, Andi Wang, Yang Li, et al. 2023. Codegeex: A pre-trained model for code generation with multilingual benchmarking on humaneval-x. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 5673–5684.