

# State of the Union of AI Coding Assistants: From Stochastic Autocomplete to Self-Evolving Agents

No Author Given

No Institute Given

**Abstract.** Automated Software Engineering (ASE) is undergoing a transformation comparable in scale to the introduction of high-level languages and early optimizing compilers. In less than five years, AI-driven coding systems have evolved from stateless autocomplete running over a single editor buffer to cloud-native software engineering agents and AI-native integrated development environments (IDEs) that operate across entire repositories, CI/CD pipelines, and production governance boundaries. This paper provides a consolidated “State of the Union” for AI coding assistants as of late 2025. We additionally emphasize the learning substrate behind modern assistants: parameter-efficient role specialization (LoRA-family adapters), long-context adaptation, and cache-aware specialization mechanisms that couple training choices to inference orchestration.

We propose a four-part taxonomy distinguishing *model-centric* systems, *repository-level* architectures, *AI-native IDEs*, and *agentic* frameworks. We analyze the architectural mechanisms that underpin these systems, including inference-time reasoning, graph-based repository maps, semantic retrieval loops, shadow workspaces, agent-computer interfaces, and self-evolving scaffolds. We connect these design choices to contemporary benchmarks such as SWE-bench and SWE-PolyBench and clarify the relationship between official leaderboards and research-oriented agent scaffolds. Finally, we examine safety, security, and governance issues—including guardrail architectures such as LlamaFirewall—and outline a research roadmap for the next generation of AI-assisted software engineering.

**Keywords:** Large language models, coding assistants, software engineering agents, AI-native IDEs, repository-level retrieval, SWE-bench, SWE-PolyBench, cloud computing, safety, guardrails.

## 1 Introduction

Automated Software Engineering (ASE) has historically advanced in waves: the adoption of high-level languages decoupled program logic from machine code; optimizing compilers introduced automatic performance tuning; refactoring tools and static analyzers automated mechanical transformations and defect detection. The most recent wave, driven by large

language models (LLMs), alters not only the mechanics of code production but also the locus of decision-making. Developers increasingly specify goals, constraints, and acceptance criteria, while AI systems propose, validate, and revise concrete implementations. LLM-based code generation has been extensively surveyed, but systematic reviews of end-to-end coding agents and IDE-integrated assistants remain lacking. To our knowledge, no prior work has unified the full spectrum of AI coding assistants across models, IDEs, and agentic systems. Unlike prior surveys that focus on code LLMs or prompt-based tools in isolation, we treat coding assistants as learning-driven systems whose behavior emerges from the interaction between (i) PEFT/long-context training choices and (ii) inference-time orchestration (routing, retrieval, tool use, cache reuse).

The first generation of LLM-powered coding tools, exemplified by the original GitHub Copilot built on OpenAI’s Codex, operated primarily as “stochastic autocomplete.” They projected local token sequences conditioned on the current buffer and a small window of surrounding code. These systems substantially improved productivity for boilerplate and idiomatic patterns but remained stateless: they did not track repository-wide invariants, understand build systems, or directly engage with test suites.

By late 2025, the landscape had bifurcated into two converging trajectories. On one side are *agentic* systems such as SWE-agent and Live-SWE-agent that treat software engineering as a sequential decision process with explicit tool use and environment interaction, evaluated on benchmarks like SWE-bench and SWE-bench Verified [1,2,3]. On the other side are *AI-native IDEs* such as Cursor, GitHub Copilot Workspace, and Codeium Windsurf [12,13,14], which embed LLMs deeply into navigation, refactoring, testing, and CI/CD workflows. These systems have moved beyond “completion plugins” to become orchestration layers for entire development lifecycles.

This paper makes three contributions:

1. It introduces a four-part taxonomy of AI coding systems, separating model-centric, repository-level, AI-native IDE, and agentic architectures and clarifying how responsibilities are divided between the “brain,” “navigator,” “workbench,” and “operator” layers.
2. It analyzes the key architectural forces shaping these systems, including static maps versus dynamic retrieval, compact models versus reasoning-heavy frontier models, and fixed versus self-evolving action spaces.

3. It synthesizes recent developments in evaluation and safety, using benchmarks such as SWE-bench and SWE-PolyBench and system-level guardrails such as LlamaFirewall to highlight where current systems succeed, where they fail, and where research is most urgently needed.

## 1.1 Scope and Limitations

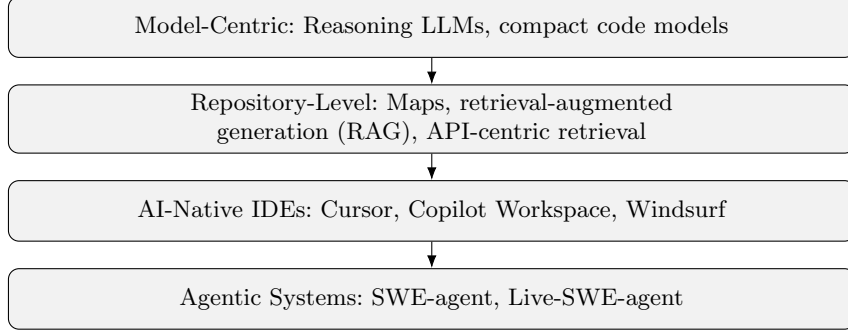
*Scope.* We survey AI coding assistants that operate over real repositories and software workflows: multi-file edits, repository-scale context construction (maps and retrieval), tool-augmented debugging/testing loops, and IDE- or CI/CD-integrated orchestration. We emphasize systems with explicit interfaces to search, file operations, execution feedback, and patch application, prioritizing publicly documented architectures and peer-reviewed artifacts from 2023–2025.

*Limitations.* The landscape evolves rapidly and results are sensitive to evaluation harnesses, repository snapshots, and scaffold implementations. Many systems report aggregate solve rates without fully specifying tool APIs, observation formats, or verification protocols, complicating reproduction and attribution. Where canonical papers are unavailable, we cite product documentation sparingly and treat it as descriptive rather than definitive.

## 2 A Taxonomy of AI Coding Systems

The phrase “AI coding assistant” now encompasses a wide range of systems with substantially different capabilities and interaction models. To make sense of this diversity, we introduce a four-layer taxonomy (Figure 1):

- *Model-centric systems* focus on the reasoning capabilities of the underlying language model.
- *Repository-level systems* manage retrieval and representation of code context at project scale.
- *AI-native IDEs* integrate models into the development environment and tooling.
- *Agentic systems* expose tool interfaces and support long-horizon planning.



**Fig. 1.** Four-layer taxonomy of AI coding systems.

These layers are conceptual; practical systems often span multiple layers. For example, an AI-native IDE may embed both a compact in-IDE model and a remote frontier model, while also providing agentic workflows for issue-to-PR automation.

## 2.1 Model-Centric Systems

Model-centric approaches emphasize improvements inside the language model. Frontier “reasoning” models allocate explicit computation to planning, often via multi-step inference or hidden chain-of-thought decoding. On SWE-bench [1], for example, models must infer which files are relevant, reason about failing tests, and synthesize patches that preserve invariants across multiple modules. These demands have driven research into inference-time scaling and “System 2”-style deliberation.

At the same time, compact models specialized for code have gained prominence. JetBrains’ Mellum-4b-base [8] and related work [9] show that a 4 billion parameter model, trained on large-scale code corpora with multi-file objectives, can provide competitive in-IDE performance for completion and navigation tasks. Mellum is architected for low-latency fill-in-the-middle (FIM) and integrates tightly with JetBrains IDEs, leveraging project indices and type information. This illustrates that data curation and integration with tooling can compensate for smaller parameter counts in certain usage regimes.

## 2.2 Repository-Level Systems

Model-centric improvements are limited by the information presented at inference time. Repository-level systems address the problem of extracting, ranking, and structuring relevant code from large, evolving codebases.

Graph-based repository maps, as implemented in *aider* [4], construct abstract syntax trees (ASTs) using Tree-sitter, extract symbols, and build dependency graphs across files. When a developer asks a question about a function or file, *aider* ranks neighboring nodes and constructs a concise “skeleton” view of the relevant files, including signatures and docstrings but omitting implementation details. This approach prioritizes structural correctness and dramatically reduces hallucinated imports and missing methods.

Dynamic retrieval, exemplified by *RepoCoder* [5], adopts a complementary strategy based on semantic similarity. *RepoCoder* first generates a tentative completion using limited context. It then uses the tentative output as a query into a vector index, retrieves semantically similar code fragments, and regenerates a refined answer with the retrieved context. Even if the initial attempt is wrong, it typically contains identifiers such as function names or variable names that unlock relevant code in the index.

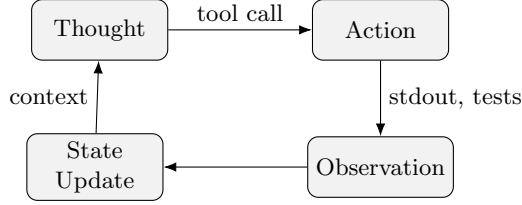
API-centric systems such as *API-Repo* [6] narrow the retrieval focus to internal APIs and their call sites. *API-Repo* indexes all function signatures and method definitions, then retrieves example usage patterns given a particular cursor location. This design acknowledges that many practical errors stem from incorrect API usage rather than misunderstandings of local control flow.

### 2.3 AI-Native IDEs

AI-native IDEs integrate language models with editor, build, and testing infrastructure in a first-class way. *Cursor* [12] demonstrates this paradigm by maintaining a shadow workspace. When the model proposes edits, *Cursor* applies them to the shadow workspace, invokes language servers, linters, and sometimes test commands, and only then surfaces validated suggestions. This feedback loop filters out syntactically invalid or type-inconsistent edits before they reach the user.

*GitHub Copilot Workspace* [14] supports more extensive workflows. Given a *GitHub* issue, *Workspace* formulates a plan, identifies affected files, generates patches, runs tests in a cloud environment, and prepares a pull request. Rather than focusing on token-level completion, it operates at the granularity of tasks and patches. *Codeium Windsurf* [13] goes further in modeling developer intent by tracking cursor movement, search queries, and edit sequences to infer high-level goals.

AI-native IDEs differ primarily in how deeply they integrate validation and workflow automation. *Cursor* emphasizes shadow-workspace valida-



**Fig. 2.** Simplified Thought–Action–Observation loop in an agentic coding system.

tion for inline edits, Copilot Workspace emphasizes issue-to-PR execution, and Windsurf emphasizes agentic flows over temporal developer context, while JetBrains + Mellum targets low-latency in-IDE completion with strong project indexing.

## 2.4 Agentic Architectures

Agentic systems treat software engineering tasks as sequential decision processes. The Agent–Computer Interface (ACI) used in SWE-agent and Live-SWE-agent [3,2] exposes a curated set of tools—for example, to open files, search for strings, run unit tests, and apply patches. Tool outputs are formatted to be token-efficient and machine-readable, enabling the model to reason over structured observations rather than raw terminal streams.

Live-SWE-agent [3] extends this paradigm by allowing the agent to synthesize new tools and modify its own scaffold during task execution. For instance, if existing search tools are insufficient for a particular project structure, the agent can write and invoke a custom Python script to perform targeted analysis and reuse that tool in subsequent steps. This self-evolving capability plays a key role in its reported 77.4% solve rate on SWE-bench Verified.

Figure 2 illustrates a typical Thought–Action–Observation loop in an agentic coding system.

## 2.5 From Scaffolded Agents to Learned Controllers

Agentic coding systems span a spectrum from mostly scaffolded pipelines to increasingly learned control. At one end are *scripted* agents: fixed Thought–Action–Observation loops over a predefined tool API, where most structure (tools, prompts, termination, and verification) is engineered. At the other end are *self-evolving* systems that can extend their

toolset or modify execution structure during deployment, effectively performing adaptation in the loop.

*What is scaffolded vs. what is learned?* In practice, end-to-end performance is jointly determined by (i) the *learned policy* inside the model (tool choice, hypothesis revision, patch synthesis) and (ii) the *scaffold* (action space, observation format, summarization/truncation rules, and verification protocol). Seemingly small scaffold choices—how error logs are serialized, how file contents are compressed, how search results are ranked, when tests are run—can dominate outcomes even when the base model is unchanged.

*Controllers and training signals.* Between model and scaffold sits an intermediate control layer: routers, controllers, and verifiers. Some are heuristic (prompt-driven “when to retrieve” / “when to test”), while others are learned via trajectory supervision. For learning, flat instruction tuning is often insufficient for long-horizon software tasks because feedback is delayed and sparse. More reliable supervision comes from software artifacts produced during execution—plans, diffs, stack traces, and test outcomes—which can be used as structured intermediate targets or as phase-conditioned signals for role adapters.

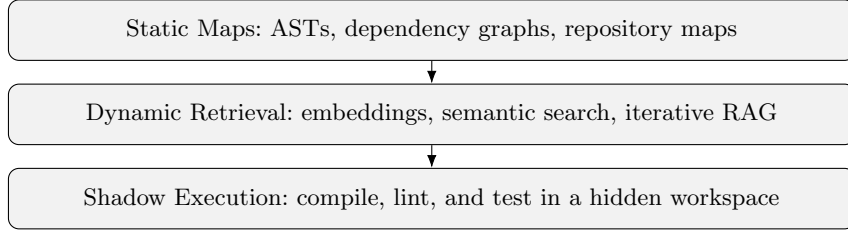
*Implications for evaluation.* Benchmark results should be interpreted as system-level properties. A high solve rate may reflect better learning, better scaffolding, or both. Reporting progress therefore benefits from ablations that separate (i) training choices (PEFT, long-context adaptation, trajectory data) from (ii) scaffold choices (tool interface, observation compression, verification loop), especially for MLCL-style claims about learning mechanisms.

### 3 Architectural Forces and Trade-offs

Although the systems described above occupy different layers, they are shaped by a common set of architectural forces. We highlight four: context management, multi-file editing, model and tool routing, and the cost-quality frontier.

#### 3.1 Context Management: Static Maps vs. Dynamic Retrieval vs. Shadow Execution

Figure 3 contrasts static maps, dynamic retrieval, and shadow execution. Static maps provide strong guarantees that important interfaces



**Fig. 3.** Complementary context strategies: static maps, dynamic retrieval, and shadow execution.

and types are visible to the model. Dynamic retrieval excels at discovering patterns and idioms spread across the repository. Shadow execution provides runtime feedback that can catch subtle type or integration errors even when the retrieved context is incomplete.

In practice, successful systems increasingly combine all three. For example, an AI-native IDE might use graph-based maps to ensure interface correctness, semantic retrieval to provide examples of similar patterns, and shadow execution to validate suggestions against real toolchains.

### 3.2 Long-Context Learning and Cache-Aware Specialization

Repository-scale coding assistants must reason over large amounts of state: repository snapshots, symbol maps, retrieved exemplars, diffs, stack traces, test logs, and long agent histories. This creates sustained pressure on context length that is qualitatively different from single-shot code generation. In long-horizon workflows, the context is not only longer; it is also *sticky*: large prefixes (repository context, prior actions, and tool outputs) persist across many turns, while only a small suffix changes from step to step.

*Window extension under PEFT.* Parameter-efficient long-context methods such as LongLoRA and LongQLoRA extend usable context windows without retraining full model weights, making them attractive for code assistants where fine-tuning budgets and deployment constraints are tight. Conceptually, these approaches adapt attention/parameterization so that models trained on shorter contexts remain stable over longer sequences, enabling the assistant to condition on project-wide structure and longer execution traces. In practice, however, simply *making the window larger* does not guarantee better agent behavior: the agent still must decide what to attend to, what to compress, and how to preserve invariants across multi-step edits.



*Why window extension is insufficient.* Long-horizon software tasks stress not just capacity but *reuse*. Many agent steps share a common prefix (issue description, repository map, previously opened files, prior diffs) while differing only in incremental observations (new test failure, a specific file snippet, the latest patch). If the model re-encodes the entire prompt at every step, latency and cost become dominated by repeated prefill, and the system becomes brittle to prompt bloat. Moreover, larger contexts can amplify distraction: irrelevant retrieved content and stale traces may compete with the true causal signals (the failing assertion, the impacted function, the minimal patch set). This is why long-context modeling for coding assistants should be treated as a systems-learning co-design problem, not a pure architecture change.

*Cache-aware specialization (aLoRA) and prefix reuse.* Cache-aware specialization mechanisms such as Activated LoRA (aLoRA) directly target the sticky-prefix regime. The core idea is to decouple *context encoding* from *behavioral adaptation*: encode the shared prefix once, reuse the base model’s key-value cache across steps or roles, and activate adapters only for newly generated tokens (or after an activation boundary). This preserves prefix KV reuse while still allowing phase- or role-conditioned behavior (planner vs. editor vs. debugger). From a learning perspective, this separation biases adapters toward decision-making and tool-use policies rather than repeatedly re-encoding static state; from a systems perspective, it reduces the dominant prefill cost in long-horizon agent loops and makes mono-host deployments more attractive as context windows grow. The broader takeaway is that future long-context gains for coding agents will likely come from *where* and *when* adaptation is applied (activation boundaries, learned gating, observation compression), as much as from increasing the nominal context length.

### 3.3 Multi-File Editing and Control Surfaces

Multi-file refactoring introduces both opportunity and risk. Scoped refactoring, where the user explicitly selects a small set of files, aligns with existing review practices and minimizes risk, but may miss necessary changes in other modules. Global refactoring can perform large-scale migrations—for example, renaming a core type across hundreds of files—but carries the risk of overgeneralization.

Systems such as Copilot Workspace constrain refactoring through plans and human review checkpoints [14]. Windsurf and Cursor offer project-wide operations but surface detailed diff previews and sometimes

group changes into separately reviewable chunks [12,13]. These designs recognize that control surfaces and user experience are as important as model quality in safely deploying powerful refactoring capabilities.

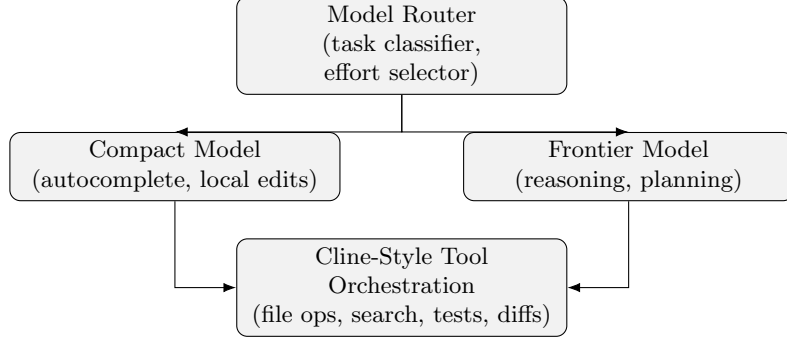
### 3.4 Model Routing and Cline-Style Tool Orchestration

Modern AI coding systems increasingly decouple the notion of a single model from the execution of a task. Rather than relying on one frontier model, systems such as Cursor, Claude Workbench, and Windsurf employ *model routers*: lightweight classifiers or heuristics that determine which model, which temperature, and which reasoning strategy should be invoked for a given step. This emerging paradigm treats coding not as a monolithic language-model problem but as an orchestration problem in which multiple specialized models collaborate.

A typical routing layer performs three functions. First, it categorizes the request into a coarse task type—such as autocomplete, documentation synthesis, multi-file refactoring, or issue-to-PR planning. Second, it selects an appropriate inference strategy. Low-latency tasks are delegated to compact models (e.g., Mellum or distilled 1–3B variants), whereas reasoning-intensive tasks may invoke large frontier models with multi-pass, chain-of-thought, or “effort-controlled” settings [16,17]. Finally, the router chooses a tool-execution path: whether to run static analysis, invoke the test harness, call a code search primitive, or construct a more complex toolchain.

Recent systems extend routing beyond model selection to full *tool orchestration*. Claude Workbench, for example, introduces a structured function-calling pipeline (often referred to informally as *cline*) in which the model emits a sequence of typed tool invocations—file edits, code searches, directory scans, test execution, or semantic diffs—as part of a single episode [17]. Each tool produces a machine-readable observation, which is then fed back into the model to inform the next step. This transforms the model from a single-shot generator into a stateful controller operating within a typed software environment.

From an architectural perspective, Cline-style orchestration occupies a middle ground between IDE-native integration and fully agentic scaffolds. It provides the structure and safety guarantees of an IDE while granting the autonomy and flexibility associated with research agents such as SWE-agent and Live-SWE-agent. Early evidence suggests that routing and orchestration are responsible for a substantial portion of the performance gains seen in 2025–2026 systems, particularly for long-horizon



**Fig. 4.** Inference-time model routing and Cline-style tool orchestration.

tasks where the cost of invoking large models repeatedly can dominate total latency and cloud cost.

### 3.5 Parameter-Efficient Role Specialization

Beyond routing across distinct models, an increasingly common design is to specialize a single backbone model into multiple roles using parameter-efficient fine-tuning (PEFT). Low-Rank Adaptation (LoRA) introduces small, trainable rank-decomposed updates while freezing the base model weights, enabling efficient role conditioning without duplicating full model parameters.

In the context of coding assistants, LoRA-style specialization supports planner-, navigator-, and operator-like behaviors within a shared representation space. This is particularly relevant for agentic workflows in which role switching occurs frequently and maintaining consistency across steps is critical. Systems such as MapCoder-Lite demonstrate that multi-agent decomposition can be preserved even under strict memory and latency budgets by attaching role-specific adapters to a single compact model. A central learning challenge is preventing role adapters from drifting into inconsistent patch styles or tool-use behaviors across phases, suggesting supervision from structured trajectories (plans, diffs, test traces) rather than flat instruction tuning alone.

From a learning perspective, role-conditioned adapters decouple task specialization from general code understanding. The base model captures broad syntactic and semantic structure, while adapters encode task- or phase-specific inductive biases.

### 3.6 Learning Dynamics of Role-Conditioned Adapters

While parameter-efficient fine-tuning is often motivated by memory and deployment efficiency, its impact on *learning dynamics* is particularly pronounced in agentic coding systems. Role-conditioned adapters do not merely specialize outputs; they implicitly shape the model’s internal control flow by biasing attention, activation patterns, and token-level decision boundaries toward phase-specific behaviors such as planning, editing, or debugging.

In staged coding pipelines such as MapCoder and CodeSim, different phases exhibit qualitatively different loss landscapes. Planning emphasizes abstraction, dependency selection, and long-range coherence, while editing prioritizes syntactic precision and local correctness. Debugging further shifts the distribution toward error localization and hypothesis testing. Training a single model to satisfy all phases via flat instruction tuning often leads to interference, where improvements in one phase degrade performance in others. Role-conditioned adapters provide a mechanism to isolate these inductive biases without fragmenting the representation space across multiple full models.

A central learning challenge is preventing role adapters from drifting into inconsistent patch styles or tool-use behaviors across phases. Empirically, such drift manifests as planner adapters that overfit to concrete edits, or editor adapters that hallucinate global refactors. This suggests that supervision derived from *structured trajectories* —including explicit plans, intermediate diffs, and test execution traces—is more effective than flat instruction-response pairs. Trajectory-level supervision preserves the temporal and causal structure of agent behavior, enabling adapters to learn not only what to produce but *when* and *why* to act.

From a systems perspective, role-conditioned adapters also reduce the need for explicit external controllers. Instead of routing requests across heterogeneous model endpoints, a mono-host system can learn soft role transitions through adapter gating mechanisms, allowing the backbone model to maintain shared context while dynamically modulating behavior. This tight coupling between learning and inference orchestration is difficult to achieve in multi-host designs where roles are enforced architecturally rather than learned.

### 3.7 Mono-Host vs. Multi-Host Architectures

Agentic coding systems can be implemented either as collections of independent model endpoints (multi-host) or as mono-host systems that

retain a single backbone model and switch behavior dynamically. Multi-host designs simplify isolation and independent scaling but often incur repeated prompt prefill and redundant memory usage when long shared context must be re-encoded across roles.

Mono-host architectures emphasize reuse. A shared backbone processes the common repository context once, while lightweight specialization mechanisms modulate behavior for planning, editing, testing, or review. This design becomes increasingly attractive as context windows grow and prefill cost dominates latency.

The trade-off is architectural rather than conceptual. Multi-host systems align naturally with microservice deployment and fault isolation, whereas mono-host systems align with learning efficiency, cache reuse, and tight integration between reasoning phases. Hybrid designs—mono-host specialization within a node and multi-host scaling across nodes—are likely to dominate production systems.

This choice also affects learning: mono-host designs favor adapter banks and learned gating over full model routing, while multi-host designs favor independently fine-tuned endpoints and explicit controller policies.

### **3.8 Cost, Latency, and Reasoning Depth**

Agentic workflows are computationally intensive. They may require dozens of tool invocations, multiple retrieval cycles, and many hundreds or thousands of tokens of internal reasoning. Routing strategies increasingly determine which model to invoke at which step, trading off latency against expected difficulty.

Recent model families such as Gemini 2.5 and Gemini 3 [15,16] and Claude 4.5 [17] explicitly expose “effort” or “reasoning” controls, allowing systems to spend more compute on difficult tasks such as multi-file bug fixes while using lighter models for simple completions. The design of these routing policies remains an open area at the intersection of systems and machine learning.

## **4 Evaluation and Benchmarking**

### **4.1 SWE-bench and SWE-bench Verified**

SWE-bench [1] consists of 2,294 real-world GitHub issues drawn from popular Python repositories such as Django, scikit-learn, and Matplotlib. Each task provides a repository snapshot, an issue description, and a hidden test suite. An agent “solves” a task if the patch it proposes compiles

and passes the hidden tests. The SWE-bench Verified subset [2] narrows this to 500 tasks that are manually vetted for clarity and solvability.

Over time, reported performance on SWE-bench Verified has increased substantially. Early works achieved below 30% solve rate under strict evaluation conditions. Live-SWE-agent reports 77.4% solve rate using a self-evolving scaffold paired with a frontier model [3]. Vendor-reported performance for Claude 4.5 and Gemini 3 on similar harnesses indicates that model improvements alone also contribute meaningfully to progress [16,17].

## 4.2 SWE-bench Pro

SWE-bench Pro shifts focus from isolated bug fixes to more complex, long-horizon tasks, such as implementing new features or modifying API behaviors across multiple modules. Early experiments show significant drops in performance compared to SWE-bench Verified, suggesting that long-range planning, global reasoning, and cross-module invariants remain unsolved challenges.

## 4.3 SWE-PolyBench

SWE-PolyBench [7] extends evaluation to multilingual, repository-level settings. It comprises more than two thousand tasks across languages including Java, JavaScript, TypeScript, and Python, and introduces Concrete Syntax Tree (CST) node-level retrieval metrics. Rather than checking only whether an agent opened the correct file, CST-level metrics evaluate whether it located the specific function or class relevant to the fix.

Agents that excel at Python often fail dramatically on Java and TypeScript tasks, even when the same model is used. Compilation barriers and stricter type systems reduce the amount of runtime feedback available to guide iterative correction. Table 2 summarizes high-level characteristics of SWE-bench, SWE-bench Pro, and SWE-PolyBench.

## 5 Safety, Security, and Governance

As coding assistants gain write access to critical repositories and infrastructure, safety and governance concerns become central rather than peripheral. Misaligned or brittle systems can introduce security vulnerabilities, propagate subtle bugs, or leak sensitive information.

Model-level alignment aims to discourage harmful outputs, but several studies have observed a “safety tax” in which strongly aligned models

**Table 1.** High-level characteristics of key coding benchmarks.

Benchmark	Langs	Tasks	Primary Focus
SWE-bench	Python	2,294	Issue-based bug fixing
Verified	Python	500	Curated solvable subset
Pro	Python	<1,000	Long-horizon, multi-file changes
SWE-PolyBench	4	2,000+	Multilingual, repo-level evaluation

exhibit degraded reasoning performance on complex technical tasks. For software engineering, this can manifest as incorrect bug fixes, missed edge cases, or brittle behavior under adversarial inputs. This tension motivates system-level guardrails that complement model-level alignment.

LlamaFirewall [10,11] represents one such guardrail architecture. It provides input filters for prompt injection and jailbreak attempts, monitors intermediate traces for signs of goal misalignment, and applies static analysis to generated code to detect insecure patterns such as hard-coded secrets or unsafe deserialization. Critically, LlamaFirewall operates outside the model, allowing safety policies to be updated or audited independently of the base LLM.

Privacy and governance constraints further shape system design. Cursor publishes clear documentation on data usage and offers modes in which user code is not stored or used for training [12]. JetBrains offers Mellum deployments that run locally or in private clouds [8,9], satisfying requirements for strict data locality. Agent observability frameworks such as Phoenix and LangSmith [18,19] provide trace-level insights into tool use and decision-making, enabling organizations to investigate and remediate unexpected behavior.

## 6 Related Work

The evolution of AI coding systems builds upon several strands of prior work. Early statistical language models for code, including n-gram models and neural sequence models, established the feasibility of learning token distributions over source code. Program synthesis research explored constrained decoding guided by specifications or examples, but was constrained by scale and domain specificity.

The advent of transformer-based language models dramatically extended the scope of code generation, enabling general-purpose models that could perform completion, translation, summarization, and synthesis for many languages. Early code-focused transformers, however, still operated primarily on single files or small snippets. Repository-scale systems such as RepoCoder [5] and aider [4] represent the next step in bringing whole-project context into the loop.

From the perspective of software engineering research, AI-native IDEs connect to long-standing work on developer productivity, human–computer interaction, and cognitive load. Studies of code navigation, comprehension, and refactoring have long advocated for tool designs that surface relevant context at the right time. AI-native IDEs extend this by offloading not only retrieval but also synthesis and validation to models.

Agentic architectures, finally, intersect with research on automated debugging, test generation, and intelligent tutoring systems. Many of these systems historically relied on symbolic reasoning or carefully crafted heuristics. Modern LLM-based agents replace or augment symbolic components with learned policies, but retain the core idea that software engineering tasks can be modeled as multi-step processes involving exploration, hypothesis formation, and validation.

A related line studies multi-agent coding as an explicit staged process. MapCoder decomposes solving into retrieval, planning, coding, and debugging, while CodeSim strengthens the pipeline via simulation-driven plan checking and stepwise debugging. MapCoder-Lite further shows that such staged decomposition can be made parameter-efficient using role-specific LoRA adapters.

## 7 Discussion and Research Roadmap (2026–2030)

The convergence of model-centric reasoning, repository-level retrieval, AI-native IDE orchestration, and agentic scaffolds suggests several research directions for the next five years.

First, multi-model and multi-agent architectures are likely to become standard. Compact local models can handle low-latency completion, while larger frontier models and specialized analysis tools can be invoked for complex changes or security-critical reviews. Designing effective routing policies and interaction protocols between these components is an open problem at the intersection of systems, economics, and ML.

Second, context management must evolve to support not only static structure but also dynamic behavior. Hybrid systems that combine AST-



and CST-based maps, semantic retrieval, and runtime traces may better capture the invariants that matter for correctness. Integrating compiler and type-checker feedback in ways that are optimized for LLM consumption—for example, by generating machine-readable error explanations—could mitigate the language gap exposed by SWE-PolyBench.

Third, safety and governance will require deeper integration between technical and organizational layers. Guardrail systems must move beyond pattern-based blocking to incorporate policy, provenance, and accountability. Observability tooling will need to scale to large agent fleets and provide abstractions suitable for engineers, security teams, and regulators alike.

Finally, there is an opportunity to reconceptualize software engineering education and practice. As AI coding assistants become standard, curricula may shift from teaching syntax and APIs toward emphasizing specification, verification, and system-level reasoning. Research on how humans and AI systems co-develop and maintain software will be essential to ensuring that AI-augmented development leads to more robust, understandable, and secure systems.

## 8 Conclusion

This paper presents a consolidated, system-level view of AI coding assistants at a moment when the field is undergoing a structural transition. What began as stochastic, token-level autocomplete has evolved into a family of learning-driven systems that reason over repositories, interact with execution environments, and coordinate multi-phase problem-solving workflows. Despite rapid progress, the literature remains fragmented: model papers focus on benchmarks, systems papers focus on tooling, and agent papers focus on execution scaffolds. This work is the first to unify these perspectives into a coherent architectural and learning-centric landscape.

A central thesis of this survey is that modern coding assistants cannot be understood as monolithic language models. Their behavior emerges from the interaction of at least four layers: the base model, repository-scale context mechanisms, orchestration infrastructure (IDEs and toolchains), and agentic control loops. Improvements at any single layer—larger models, longer context windows, or more tools—do not translate reliably into better end-to-end performance without corresponding advances in learning objectives and systems integration.

From a learning perspective, the field is shifting away from static instruction tuning toward structured, trajectory-aware supervision. Agentic coding systems expose failure modes that are invisible in single-turn benchmarks: inconsistent patch styles, unstable tool use, compounding planning errors, and brittle cross-phase assumptions. Parameter-efficient role specialization, particularly LoRA-family adapters, offers a scalable mechanism for decomposing behavior across phases, but only when paired with learning signals that enforce cross-role coherence. Flat prompt-response datasets are insufficient for this setting; plans, diffs, test traces, and execution feedback must be treated as first-class learning artifacts.

Long-context learning further reinforces the need for systems-learning co-design. While methods such as LongLoRA and LongQLoRA demonstrate that context windows can be extended efficiently, their practical value depends on inference orchestration. Cache-aware specialization mechanisms such as Activated LoRA highlight a broader principle: as context length grows, learning should increasingly bias toward decision-making rather than re-encoding shared state. This observation reframes long-context modeling not as a pure capacity problem, but as a representation and control problem—one that aligns naturally with agentic workflows.

At the architectural level, the survey clarifies distinctions that are often blurred in benchmark-driven discussions. Mono-host and multi-host designs represent fundamentally different trade-offs in learning efficiency, cache reuse, fault isolation, and deployment complexity. Similarly, agentic systems range from scripted pipelines to self-evolving controllers, each imposing different constraints on learning stability and evaluation. Making these distinctions explicit is essential for interpreting reported performance and for designing reproducible, comparable systems.

Evaluation remains a limiting factor. Benchmarks such as SWE-bench and SWE-PolyBench have played a crucial role in advancing the field, yet they capture only a subset of the capabilities required in real-world software engineering. Long-horizon planning, cross-repository consistency, and interaction with evolving toolchains remain weakly measured. Progress will require evaluation protocols that align more closely with the learning problems surfaced by agentic systems, including robustness, consistency, and error recovery across phases.

Looking forward, we argue that the next phase of research on AI coding assistants will be defined less by raw model scaling and more by advances in controllable learning, modular adaptation, and principled orchestration. Key open directions include: learning stable role decom-

positions from trajectories, integrating compiler and runtime feedback into training objectives, designing cache-aware adaptation mechanisms for long-horizon agents, and developing evaluation frameworks that reflect system-level behavior rather than isolated task success.

By articulating a unified taxonomy, identifying shared architectural forces, and grounding the discussion in learning dynamics rather than product capabilities, this paper aims to establish a common vocabulary for the field. We view this work not as a conclusion, but as a foundation: a reference point for future research that treats AI coding assistants as what they have become—complex, learning-driven systems operating at the intersection of machine learning, programming languages, and distributed systems.

## References

1. C. E. Jimenez, J. Yang, A. Wettig, S. Yao, K. Pei, O. Press, and K. R. Narasimhan, “SWE-bench: Can Language Models Resolve Real-World GitHub Issues?,” in *Proc. ICLR*, 2024.
2. SWE-bench Team, “SWE-bench: Can Language Models Resolve Real-World GitHub Issues?,” official website, accessed Nov. 2025. [Online]. Available: <https://www.swebench.com/>
3. C. S. Xia *et al.*, “Live-SWE-agent: Can Software Engineering Agents Self-Evolve on the Fly?,” *arXiv:2511.13646*, 2025.
4. Aider Project, “Repository Maps with Tree-sitter,” project documentation, accessed Nov. 2025. [Online]. Available: <https://aider.chat/docs/>
5. F. Zhang, B. Chen, Y. Zhang, J. Keung, J. Liu, D. Zan, Y. Mao, J. G. Lou, and W. Chen, “RepoCoder: Repository-Level Code Completion Through Iterative Retrieval and Generation,” in *Proc. EMNLP*, 2023.
6. Z. Li *et al.*, “API-Repo: API-centric Repository-level Code Completion,” in *Proc. Internetware*, 2025.
7. M. S. Rashid *et al.*, “SWE-PolyBench: A Multi-Language Benchmark for Repository-Level Evaluation of Coding Agents,” *arXiv:2504.08703*, 2025.
8. JetBrains, “Mellum-4b-base: Open-Source Code Model,” Hugging Face Model Card, 2025. [Online]. Available: <https://huggingface.co/JetBrains/Mellum-4b-base>
9. N. Pavlichenko *et al.*, “Mellum: Production-Grade In-IDE Contextual Code Completion with Multi-File Project Understanding,” *arXiv preprint*, 2025.
10. S. Chennabasappa *et al.*, “LlamaFirewall: An Open Source Guardrail System for Building Secure AI Agents,” *arXiv:2505.03574*, 2025.
11. Meta AI, “LlamaFirewall: About and Documentation,” 2025. [Online]. Available: <https://meta-llama.github.io/PurpleLlama/LlamaFirewall/>
12. Anysphere, “Cursor: Data Use and Privacy Overview,” Oct. 2025. [Online]. Available: <https://cursor.com/data-use>
13. Codeium, “Windsurf Editor: The First Agentic IDE,” product site and blog, 2024–2025. [Online]. Available: <https://windsurf.com/>
14. C. Reddington and C. Bemis, “5 Tips and Tricks When Using GitHub Copilot Workspace,” GitHub Blog, Oct. 2024.

15. Google DeepMind, “Gemini 2.5: Our Most Intelligent AI Model,” Google Blog, Mar. 2025.
16. Google DeepMind, “A New Era of Intelligence with Gemini 3,” Google Blog, Nov. 2025.
17. Anthropic, “What’s New in Claude 4.5,” Claude Models Documentation, Nov. 2025. [Online]. Available: <https://platform.claude.com/docs/>
18. Google ADK Docs, “Agent Observability with Phoenix,” Agent Development Kit Documentation, 2025. [Online]. Available: <https://google.github.io/adk-docs/observability/phoenix/>
19. LangChain, “LangSmith Observability,” LangSmith Product Docs, 2025. [Online]. Available: <https://www.langchain.com/langsmith/observability>
20. E. J. Hu *et al.*, “LoRA: Low-Rank Adaptation of Large Language Models,” in *Proc. ICLR*, 2022.
21. W. Lee, J. Cho, and J. Choi, “MapCoder-Lite: Squeezing Multi-Agent Coding into a Single Small LLM,” *arXiv:2509.17489*, 2025.
22. J. Yang, “LongQLoRA: Efficient and Effective Method to Extend Context Length of Large Language Models,” *arXiv:2311.04879*, 2023.
23. S. Yao *et al.*, “ReAct: Synergizing Reasoning and Acting in Language Models,” in *Proc. ICLR*, 2023.
24. T. Schick, J. Dwivedi-Yu, R. Dessì, R. Raileanu, M. Lomeli, L. Zettlemoyer, N. Cancedda, and T. Scialom, “Toolformer: Language Models Can Teach Themselves to Use Tools,” *arXiv:2302.04761*, 2023. :contentReference[oaicite:0]index=0
25. N. Shinn, F. Cassano, A. Labash, B. Gopinath, K. Narasimhan, and S. Yao, “Reflexion: Language Agents with Verbal Reinforcement Learning,” *arXiv:2303.11366*, 2023.
26. X. Wang *et al.*, “Self-Reflective Large Language Models,” in *Proc. NeurIPS*, 2023.
27. S. Yao *et al.*, “Tree of Thoughts: Deliberate Problem Solving with Large Language Models,” *arXiv:2305.10601*, updated 2024.
28. Y. Zhou *et al.*, “Least-to-Most Reasoning Enables Complex Problem Solving in Large Language Models,” *arXiv:2205.10625*, extended analysis 2024.
29. J. Kim *et al.*, “AgentBench: Evaluating LLMs as Agents,” in *Proc. ICLR*, 2024.
30. Z. Liu *et al.*, “SWE-agent: Agent-Computer Interfaces Enable Automated Software Engineering,” in *Proc. NeurIPS*, 2024.
31. A. Prasad *et al.*, “CodeSim: Multi-Agent Code Generation via Simulation and Verification,” *arXiv:2402.01935*, 2024.
32. M. Islam *et al.*, “MapCoder: Multi-Agent Code Generation for Complex Tasks,” *arXiv:2401.09590*, 2024.
33. T. Dettmers *et al.*, “QLoRA: Efficient Finetuning of Quantized LLMs,” in *Proc. NeurIPS*, 2023.
34. C. Pouliot *et al.*, “Composable Low-Rank Adapters for Parameter-Efficient Transfer Learning,” *arXiv:2406.01234*, 2024.
35. K. Greenewald *et al.*, “Activated LoRA: Efficient Adaptation of Large Language Models at Inference Time,” *arXiv:2504.12397*, 2025.
36. Y. Chen *et al.*, “LongLoRA: Efficient Fine-Tuning of Long-Context Large Language Models,” *arXiv:2309.12307*, 2023.
37. J. Yang *et al.*, “LongQLoRA: Efficient and Effective Context Length Extension for LLMs,” *arXiv:2311.04879*, 2023.
38. Y. Sun *et al.*, “Ring Attention with Blockwise Transformers for Efficient Long Context Modeling,” *arXiv:2403.02042*, 2024.