# The Rise of Agentic AI Coding Tools: A New Stack for Developers

AI coding assistants have rapidly evolved from simple autocomplete plugins into **agentic** tools that can read, write, and even execute code on our behalf. In this newsletter, we'll explore the current landscape of these AI coding agents – from terminal-based copilots to IDE-native assistants and code-review bots – and discuss how they work, which models power them, and what it means for developers and teams. We'll also look at the history (from GitHub Copilot's debut to today's autonomous "AI software engineers"), emerging trends (multi-agent collaboration, open protocols), and strategic takeaways for building your own AI-enhanced dev stack.

## From Autocomplete to Autonomous: A Brief Evolution

Just a few years ago, tools like GitHub Copilot introduced AI **code completion** to suggest the next line or block of code. Copilot (powered by OpenAI's Codex model) quickly proved popular, and by 2023 nearly **half of some developers' code** was being generated by AI assistants [1]. This success set the stage for a new generation of **"agentic" AI coding tools** – systems that go beyond suggestions to actively perform coding tasks. These agents can *understand entire codebases, execute commands, run tests, fix bugs, and even generate multi-step plans* to achieve a goal. In short, we've moved from **passive AI helpers** to **proactive AI collaborators** in our development workflows.

Several factors enabled this shift: larger context windows (allowing models to ingest whole repositories), specialized coding models (tuned for software tasks), and sandboxed execution environments for safe code running. The result is an explosion of new AI coding tools and frameworks. Let's break down the current landscape by category, then dive into capabilities, models, and strategies behind each.

## Terminal-Based AI Coding Agents

For developers who live in the command line, a crop of AI agents bring ChatGPT-style intelligence **directly into the terminal**. These tools act like an AI pair programmer accessible via CLI, capable of executing shell commands and editing files based on natural language prompts.

- **OpenAI Codex CLI (Terminal GPT)** – OpenAI's official CLI tool (confusingly named after the older "Codex" model) is a lightweight, open-source coding agent that you run locally [2]. It connects to GPT-4 (or other models via API) and is *"chat-driven development that understands and executes your repo"* [2]. Codex CLI can read your codebase, generate or modify files, and run tests or scripts in a sandbox. It supports modes ranging from "suggestion-only" to full autonomy. In **Suggest mode**, it only proposes changes (you approve before anything runs), while **Auto-Edit and Full-Auto modes** let it apply code edits or execute commands without confirmation [3] [4]. To mitigate risk, OpenAI built in strong security: commands run in a restricted sandbox (macOS uses Apple's Seatbelt; Linux uses Docker with no network access) [5] [6]. The tool is free (bring your own API key) and even

supports custom model backends – you can swap in local models (like via Ollama or OpenRouter) by editing a config [7] [8] . In short, Codex CLI delivers *"ChatGPT-level reasoning plus the power to actually run code, manipulate files, and iterate – all under version control"* [2] .

- **Anthropic Claude Code** – Anthropic (the creators of Claude) released their own agentic CLI in late 2024, simply invoked by the `claude` command [9] . It integrates **Claude 4 (Opus)**, one of the most advanced coding models, to provide a wide range of dev abilities. Claude Code can *edit multiple files, fix bugs across your codebase, answer questions about your architecture, run tests and linters, search git history, resolve merge conflicts, and even generate commits/PRs – all through natural language commands* [10] . Under the hood it offers dozens of tools (file operations, code search, web browsing, etc.) to fulfill your requests [11] . Notably, Claude Code *works entirely in your local environment* with no additional server – your prompts go straight to Anthropic's API (or an on-prem model host) [12] . It maintains context on your whole project structure and "understands" the entire repo as it works [13] . With direct file writes and `git` operations, this agent truly "acts" on your behalf. Anthropic open-sourced Claude Code (available via npm) and even provided an SDK for integration. You can run it programmatically (TypeScript or Python SDK) to build your own tools, or connect it to other services via the **Model Context Protocol (MCP)** [14] [15] (an open standard for LLM tools, more on that later). In short, Claude Code brings the full power of Anthropic's AI into your terminal and development workflow.

- **Google Gemini CLI** – Not to be outdone, Google unveiled *Gemini CLI* in mid-2025 as a free, open-source terminal agent [16] [17] . Gemini CLI connects to Google's new **Gemini 2.5 Pro** model, giving users an astounding **1 million-token context window** for code understanding [18] . Google positions this as a "versatile local utility" for coding, content creation, research, and more [19] . Out-of-the-box, Gemini CLI can write and debug code, manipulate files, execute shell commands, and perform "dynamic troubleshooting" in response to errors [20] . It has built-in web browsing (grounding prompts with Google Search) and supports extensions via MCP [14] . A major selling point is Google's generous free usage: with a personal Google login, you get access to Gemini 2.5 Pro with **up to 60 requests/minute and 1000 requests/day** at no charge [18] . This *"unmatched allowance"* is meant to ensure individuals *"rarely, if ever, hit a limit"* in normal use [21] [22] . For power users or teams, it can also be configured to use your own Vertex AI or AI Studio keys for higher volume or different model versions [23] . Strategically, Gemini CLI is fully open source (Apache 2.0) [24] – Google clearly hopes to foster community contributions and make Gemini a fixture in dev workflows. It's also tightly integrated with Google's IDE assistant (Gemini Code Assist), which means you can seamlessly move between VS Code and terminal with the same AI "agent" backing you [25] [26] .

- **Aider (open-source CLI)** – Before the big players entered, open-source projects like **Aider** paved the way for terminal pair-programming. Aider is a popular CLI tool that lets you chat with GPT-4 or other models to modify code in a local git repo [27] . You can ask it to implement a feature or fix a bug, and it will edit the relevant files, staging changes as git diffs you can review. It's effectively a simpler precursor to Claude/Gemini CLI – *in fact, the Aider community is credited with inspiring many later tools* [28] . Aider can run fully locally if paired with open LLMs (it recently added support for hosting models via Ollama, enabling "100% local" operation [29] ). While it lacks the deep toolset of the newer agents, it remains a lightweight and free option for those who want AI code edits in the terminal without vendor lock-in.

*Other notable mentions:* **Cline** is an open-source VS Code extension that bridges CLI and IDE (hence "CLInE"). It wraps Claude or Gemini behind the scenes and offers dual modes – a "Plan" mode to draft a solution and an "Act" mode to execute changes – with terminal integration and MCP tool support [30] . There's also **Devin** by Cognition, a closed beta product branding itself as an "AI Software Engineer." Devin runs multiple cloud agents in parallel to tackle big coding tasks, but its *steep* pricing (around **$500/month** for access [31] [32] ) and hefty claims (a $2B valuation startup) have drawn skepticism. Most individual developers will gravitate toward the more accessible agents above, but it's worth noting the ambition in this space – some are literally trying to deliver an "autonomous engineer" as a service.

## IDE-Integrated AI Assistants and Agents

While terminal agents are powerful, many developers prefer to stay within their IDE/editor. The good news: today's AI assistants are deeply integrating into editors, from VS Code to JetBrains, effectively embedding agentic capabilities where you code.

- **Cursor** – *Cursor* is an AI-native code editor (built on VS Code under the hood) that has become a favorite for many. Developed by startup AnySphere, Cursor provides an **AI pair-programming experience** with support for multiple models. Free users can use a fast, custom model ("cursor-small") or Claude Instant, while Pro users get GPT-4 access [33] . Cursor isn't just about autocompletion; it offers an in-editor chat that understands your whole project context and can apply changes. You can highlight code and give natural language instructions (e.g. "refactor this function to use async/await"), and it will directly edit the file. Under the hood, Cursor actually analyzes your entire codebase to build context [34] [35] , enabling more global reasoning than typical copilot-style tools. It supports many languages (JS, Python, TS, etc.) and even works when connected to remote dev environments over SSH [36] . In terms of capabilities, Cursor handles code generation, explanation, bug fixes, and even documentation. It's essentially trying to *"keep you in flow"* by letting you talk to your editor as if it were a collaborator [37] [38] . Strategically, AnySphere is a well-backed startup (OpenAI's Startup Fund invested early) [39] . They've attracted users at companies like Midjourney, Shopify, and more [40] . Cursor's positioning is to be the "AI-first IDE" competing with the likes of Visual Studio + Copilot. It has a freemium model (basic usage free, $20/month Pro for expanded GPT-4 limits) [33] . As a third-party, it has to rely on model APIs (OpenAI, Anthropic) – but it differentiates via a tight UX integration and features like **natural language commands** (execute an edit or action via text) and an upcoming multi-file editing improvements (which early versions struggled with [41] ). If you want an IDE where AI is not an afterthought but a core part of the UI, Cursor is a leading choice.

- **Windsurf (formerly Codeium)** – Windsurf is another AI-powered IDE making waves. It's the evolution of **Codeium**, a popular free code completion tool. The team behind Codeium launched Windsurf Editor in 2025 as *"the first AI agent-powered IDE"* [42] . The editor (available on Mac/Win/Linux) comes with an AI agent named **Cascade** at its core. Cascade is described as *"an agent that codes, fixes, and thinks 10 steps ahead"* [43] . In practical terms, Cascade can interpret high-level tasks and execute multi-step coding workflows. For example, Windsurf has a **Cascade "Write Mode"** that functions like AutoGPT – *you give a prompt for a new feature, and it will generate multiple files, run and test them, debug errors, and iterate* until the feature works [44] . At each step it asks for approval before, say, running a command or applying a major refactor [45] (you can adjust this autonomy level, including a Turbo mode that auto-executes approved steps [46] ). Cascade also has built-in web browsing/documentation parsing to bring in external context when needed [47] . A standout feature

of Windsurf is its **"Memories"** – the agent automatically notes important details about your codebase or past instructions to improve consistency [48] [49] . You can also define **AI Rules** to guide the agent (e.g. "always use React hooks" or "follow our coding styleguide") [50] . Windsurf supports **MCP plugins** too – enabling integration with tools/APIs like Figma, Slack, Stripe, databases, etc., directly from the IDE [51] [52] . Essentially, the Windsurf team is aiming to replace your IDE, your Copilot, and even some of your devOps automations with one AI-infused platform. They claim over 1M users (from the Codeium legacy) and thousands of enterprise customers already [53] . Windsurf's AI model strategy is interesting: they have developed their own **SWE-1** series models specialized for software engineering (trained from scratch for coding tasks) [54] . This in-house model (and possibly fine-tuned open models) powers their **"Supercomplete"** code completion and Cascade's reasoning. However, Windsurf also offers *"first-class support for every major model provider"* [55] – allowing you or an enterprise to plug in OpenAI, etc., if preferred. Windsurf Editor is free for individual use (with enterprise paid plans), continuing Codeium's strategy of wide adoption. For those who don't want to leave their current IDE, Windsurf provides extensions (VS Code, JetBrains) that bring in some features (primarily code completion), but the full Cascade agent experience is in their dedicated editor [56] . In summary, Windsurf is a more **holistic AI dev environment**: it combines a coding-optimized model, an autonomous agent, and integration hooks, all wrapped in an IDE. Early adopters have noted it "makes steps easier – I just type my prompt, go away for a bit, come back and there's a web preview waiting" [57] . It's an ambitious attempt to "win the AI Coding Assistant wars" [58] with an open-core approach.

- **Gemini Code Assist (VS Code)** – This is Google's entry into IDE assistants, offered as an extension for VS Code. Initially launched in 2024 as a competitor to Copilot, **Gemini Code Assist** provided code completions and a chat panel. In 2025, Google upgraded it with an *"Agent Mode"* that essentially turns it into an **autonomous coding agent within the IDE** [26] . When you enable agent mode, you can give it a prompt like "Add a feature X and unit tests," and it will *"relentlessly work on your behalf"* – writing code across files, running tests, fixing errors, even migrating code between frameworks [26] . The agent builds a **multi-step plan** to fulfill your request and can recover from failures automatically [59] . This is the same tech behind Gemini CLI, integrated in VS Code's UI [25] [60] . The **free tier** is notable: Google allows **unlimited usage** for individuals (the free Code Assist license uses Gemini 2.5 with the highest limits in the market) [21] [61] . For enterprise, they offer paid plans or the ability to use one's own Vertex AI keys. Strategically, Google's making a play for both the pro coders (terminals) and the broader dev audience (IDE users), with a unified agent platform (Gemini + MCP + extensions). If you're a VS Code user, Code Assist offers a powerful alternative to Copilot – especially with its ability to *act* on the codebase (not just suggest) under your supervision.

- **GitHub Copilot X** – Microsoft-owned GitHub has been integrating more agentic features into Copilot, albeit gradually. Copilot began as inline code completion (still its main use), but with the **Copilot Chat** beta, it gained an IDE chat that can answer questions about your code and suggest fixes. The chat can reference multiple files and explain errors, but currently it doesn't autonomously refactor code across your project – it still relies on you to confirm and apply suggestions. However, GitHub has demonstrated upcoming features like **Copilot for Pull Requests** (which auto-summarizes diffs, suggests tests, and may even propose direct code changes in PRs). They also released a separate **Copilot CLI** (not to be confused with OpenAI's Codex CLI) that helps generate shell commands and loops via natural language. All that to say, Copilot is evolving from an autocomplete into more of an "AI assistant in the IDE," but it's doing so cautiously. It uses OpenAI's models (now GPT-4 for Copilot Chat, and a Codex-derivative for completions) and remains a closed,

paid product ($10/mo individual or enterprise licensing). While Copilot currently doesn't execute code for you, its huge user base (over 1 million developers) and integration into GitHub's ecosystem mean it will surely join the agentic trend – just within the guardrails of enterprise trust and Microsoft's platform. (One early signal: Microsoft's new Visual Studio 2022 "IntelliCode" AI beta can actually *edit your codebase* in response to a natural language prompt, a very Copilot-X-like capability, showing they're testing agent actions).

- **Replit Ghostwriter** – Replit's Ghostwriter deserves mention as an IDE-native AI, especially for those building in the browser. Ghostwriter offers code completion and a chat that can answer questions about your Replit project. It also introduced a "Generate Project from Prompt" feature where you describe an app and it scaffolds the project (a bit like GPT-Engineer). While Ghostwriter doesn't fully autonomously iterate (Replit tends to keep the human in the loop), it has a unique advantage: since Replit can run code in the cloud container, Ghostwriter can actually execute and verify code as part of the development loop. Replit hasn't marketed Ghostwriter as an agentic CLI or multi-step planner, but given their focus on beginner-friendliness, they're likely exploring ways to automate more of the coding tasks (perhaps in education contexts first). As of now Ghostwriter is a subscription service (or usage-based for free tier) using OpenAI under the hood, with some fine-tuning.

In summary, **IDE-integrated AI** is becoming table stakes. Whether via a first-party tool (Copilot, Code Assist) or a third-party (Cursor, Windsurf), developers can now get AI help *inside* their editor that not only autocompletes code but can understand project-wide context and handle non-trivial tasks. The **UX layers** vary – some prefer a separate "AI console" (Cursor's side panel or Windsurf's Cascade UI), others integrate into existing panels (VS Code's chat, inline suggestions). The key is that the AI is no longer a separate website (like one might use ChatGPT in browser) – it's in your development flow, potentially editing and running code directly with minimal friction.

## AI Code Reviewers and GitOps Agents

Another niche seeing AI disruption is **code review** – using AI to automate the checking of code changes, find bugs, and even suggest improvements in pull requests. These tools act like an "AI reviewer" that supplements your peers, catching issues early and speeding up merges.

- **CodeRabbit** – CodeRabbit is a standout platform in AI code review. It's a SaaS tool that integrates with GitHub/GitLab: when you open a pull request, CodeRabbit's agent springs into action. It uses a large-model backend (the company has worked closely with Google, hinting at using models like PaLM or Gemini via Vertex AI) to analyze not just the diff, but the entire repository context around it [62] [63] . What makes CodeRabbit "agentic" is that it doesn't only rely on the LLM's static analysis – it *actively executes tools and scripts to probe the code*. According to Google Cloud, CodeRabbit will spin up an isolated environment for each PR, *clone the repo, install dependencies, run 20+ linters and security scanners, and even execute the code if needed* [64] . The AI agent then generates specialized shell scripts on the fly to search the code for patterns (using grep, AST parsers, etc.) and extract deeper information [65] . These AI-generated scripts might, for example, detect if a change in one module impacts another, or if a new function lacks tests. CodeRabbit can even go one step further: it has the ability to **call external APIs or services** as part of review (for instance, posting a Slack message or creating a Jira ticket) by generating `curl` commands in its execution sandbox [66] . All this is done in a secure, sandboxed Cloud Run environment, since running untrusted code poses risks [67] [68] . Finally, CodeRabbit surfaces the results back in the PR – typically as an **AI-generated**

**review comment** summarizing the change, highlighting potential bugs, style issues, or missing edge cases, and even asking the author questions. Developers can then converse with the AI in the PR thread, e.g. "Explain why you think this null check is needed," and the bot will respond with context [69] [70] . CodeRabbit's value proposition is faster reviews and higher code quality: their customers cite 50% less time in review and lots of bugs caught early [71] [72] . The service has a free tier (they even made AI reviews free inside VS Code, Cursor, and Windsurf editor extensions [73] ) and paid plans for teams (with SOC2 compliance, etc. for enterprises [74] [75] ). Strategically, CodeRabbit is a startup, but by leveraging cloud scale and focusing on the pain point of PRs, it's carving a niche that even big players haven't fully tackled yet. GitHub is testing AI for PR summaries, but CodeRabbit is *far* more comprehensive in analyzing and even executing code for a thorough review [76] [77] . It essentially acts as an **automated devops engineer** that checks your work and sometimes fixes it (the team is working on "reviewers that can generate patches for issues they find").

- **Ellipsis** – Another Y-Combinator backed tool, Ellipsis, also offers AI code reviews with a focus on **bug-fixing**. It connects to your GitHub repo as a bot, and for each PR it not only comments on problems but can also push follow-up commits to fix simple issues. For example, if it detects a potential bug or style violation, Ellipsis might suggest a change – and with maintainer approval, it can commit that change. Ellipsis uses GPT-4 and emphasizes catching logical errors, security issues, and documentation mismatches (like if your code doesn't match the README, it will notice) [78] [79] . It claims to help teams merge ~13% faster by reducing the back-and-forth on minor issues [80] . The approach is similar to CodeRabbit in spirit (AI doing grunt work of review), though Ellipsis appears to lean more on static analysis + GPT reasoning rather than heavy dynamic analysis. It's also offered as a paid service with a free trial.

- **GitHub & DevOps Integrations** – Traditional platforms are starting to integrate AI reviewers as well. GitHub's upcoming **"Hey GitHub, review this PR"** feature will add a Copilot-generated analysis in pull request UI, highlighting potential issues or areas of interest. It's expected to use GPT-4 with fine-tuning on common code review comments. There's also **Amazon CodeGuru**, an older tool (non-LLM) for automated code review that provides performance and correctness tips; one could imagine Amazon infusing CodeWhisperer's AI to make CodeGuru more conversational. In CI pipelines, some teams are deploying their own AI scripts: e.g. a CI job that calls an LLM to analyze test failures or diff impacts and post results (with the new ChatGPT function calling or tools like LangChain, this is quite feasible). We're even seeing **AI ops bots** – for instance, an AI that watches your error logs in production and opens a GitHub issue with analysis when it spots a pattern (a sort of AI SRE). These remain early, but tie into the idea of multi-agent workflows (one agent writes code, another reviews it, another monitors it – all autonomously coordinating under some human oversight).

**Takeaway:** AI will increasingly assist at *every step of the code lifecycle*, not just during initial coding. Code review is a natural place since it's often repetitive ("style nit here, missing check there") and time-consuming. However, these AI reviewers are still evolving – they can miss deeper architectural issues and sometimes give false positives. Thus, they're positioned as **augmented assistants** to human reviewers, not replacements. Teams adopting them report higher consistency in reviews and less tedious feedback, freeing humans to focus on design and high-level feedback [81] [82] . As models improve, we might trust them with more automated fixes (imagine a PR that auto-merges if the AI reviewer and tests both give a green light!). Until then, using an AI reviewer is like having an extra diligent junior dev on the team who never gets tired of combing through code.

## Extensible Agent Frameworks and SDKs

Beyond out-of-the-box products, there's a growing movement to build **extensible AI agent frameworks** for coding. These are toolkits or protocols that let you customize how an AI interacts with your development environment, or even compose multiple agents together.

A few notable developments:

- **Anthropic's Claude Code SDK & MCP** – Anthropic has provided an SDK so developers can *embed Claude Code's capabilities into their own tools* [83] [84] . For example, you could build a custom GUI that uses Claude Code under the hood to perform refactoring tasks, or integrate it into your CI pipeline to auto-fix lint errors after each commit. The SDK allows running Claude Code in a non-interactive way (headless mode, as a subprocess) and controlling options like the working directory, system prompts, and even the model backend [84] [85] . This opens the door for specialized AI dev tools – think AI-assisted **migrations** (point it at a codebase and have it upgrade all usage of Library X to Library Y), or an AI **documentation generator** that reads code and produces docs continuously. Moreover, Anthropic's support of the **Model Context Protocol (MCP)** is significant [15] . MCP is an open protocol (championed by devs in this agent community) that standardizes how an LLM agent can discover and use external tools and data. With MCP, you can write a tool server (say one that provides access to a database or a design system) and any MCP-compatible AI agent (Claude Code, Gemini CLI, Cline, Windsurf's Cascade, etc.) can query it. This interoperability means you could extend multiple AI agents with the *same* custom tool and share it across your workflow.

- **OpenAI Function Calling & API** – OpenAI hasn't released an "agent SDK" per se for coding, but they have introduced **function calling** in their API which developers are using to create mini-agents. For instance, you can define functions like `read_file(path)` , `write_file(path, content)` , `run_tests()` and have GPT-4 "call" these functions during a conversation. This way, you can script an agent that, say, reads a file, gives an answer or modifies it, and writes it back – all in one prompt loop. Several open-source projects (in Python, JS) wrap GPT-4 with such function definitions to replicate what Codex/Claude Code do, albeit less robustly. OpenAI has also open-sourced the **Codex CLI** (as described), so developers can fork and tweak it. Its config system allows plugging in local LLMs (so you could experiment with an open-source model in place of GPT-4) [86] [87] . Don't be surprised if OpenAI eventually offers a more official "Dev Agent API" given the traction – for now, the community has essentially built it atop their general API.

- **LangChain, AutoGPT and Beyond** – The craze of AutoGPT in 2023 (where an AI agent spawns sub-agents to plan and solve tasks) naturally spilled into coding. Projects like **GPT-Engineer** and **Camel** tried using one agent to generate specs and another to implement code. While those early experiments were brittle, the idea of **multi-agent teams** is gaining traction. We see this in products like Devin (which touts parallel agents) and even in Windsurf's design (Cascade can utilize multiple "threads" of thinking with its tools). In open source, frameworks like **LangChain** or **GPT-4 Tools** allow savvy developers to orchestrate an agent that uses various tools (terminal, browser, etc.). For example, one could configure an agent that, given a Jira ticket, will use a GitHub API tool to create a feature branch, an editor tool to modify code, a testing tool to run tests, and a Git tool to open a PR – essentially automating an entire dev task end-to-end. These are complex to set up and error-prone today, but they hint at where things are headed.

- **Open-Source Code Models** – A crucial piece of extensibility is having models you can run anywhere. While GPT-4 and Claude are king now, the open-source community has made strides with models like **StarCoder** (15B), **Code Llama** (7B–34B) and others specialized for code. These can't yet match GPT-4's reliability, but they are improving rapidly – and they can be self-hosted. The **DevON** project (not to be confused with Devin) is an open-source AI pair programmer that emphasizes local models [88] [89] . It supports basic coding Q&A and multi-file edits for Python, with the goal of running fully offline. As 34B+ parameter code models become more viable (and as techniques like context distillation or retrieval improve their effective context window), we may see **fully local agentic IDEs** that don't require any cloud API. This is attractive for companies with strict IP policies who currently shy away from sending code to third-party APIs. Already, **Aider** and **Codex CLI** allow using local backends; Windsurf might in future (since Codeium had some self-hosted options for enterprise). Open-source AI agents plus open models could replicate 80% of what the proprietary ones do, which puts competitive pressure on the big providers to keep innovating (or open up models further).

## Strategic Postures: Open vs Closed, Startup vs Big Tech

It's fascinating to observe the strategies at play in this fast-moving domain:

- **Big Tech (OpenAI, Google, Anthropic)**: Interestingly, they have open-sourced almost all the *agent tooling* (CLI wrappers, IDE plugins) while keeping their **models** proprietary. The CLI/IDE tools are essentially funnels to increase usage of their AI models. For example, Google made Gemini CLI open and gave away huge free usage [21] [23] – they clearly want hordes of developers trying Gemini in their terminal, which could drive adoption of Google's ecosystem (and possibly yield feedback/data to improve the model). Anthropic open-sourced Claude Code and even enabled third-party model hosting (Bedrock, Vertex integration) [90] , positioning themselves as somewhat open and enterprise-friendly. OpenAI open-sourced Codex CLI likely to ensure they remain the default choice for developers even as competitors multiply. By being open with the *integration layer*, these companies allow the community to contribute and build around their models (plugins, extensions, etc.), which further entrenches their core model offerings. It's a savvy move: they relinquish a bit of control (the UI/UX layer) but keep the crown jewels (the model and its API).

- **Startups**: The startups in this space differentiate in a few ways – some build their own models (Windsurf/Codeium's SWE models, possibly CodeRabbit fine-tuning things for analysis), some focus on unique UX or domain (Cursor's "AI editor for programmers", CodeRabbit's focus on PR quality, Ellipsis on bug-fixing, etc.), and some aim for high-end "full autonomy" (Devin). A big challenge for startups is that model providers can quickly implement similar agent features, *and* they have the advantage of tighter integration (e.g. OpenAI can bake in function calling or have privileged model abilities; GitHub can embed Copilot deeply in PR workflow). To survive, startups often go **open-source or freemium** to attract users quickly. Codeium/Windsurf made their core free, gaining a million users – a base from which they upsell enterprise features. Cursor also offers a generous free tier (with some GPT-4 usage) to get mindshare. Many agent startups actively support multiple models (to avoid dependence on a single API and to appeal to users' preferences/cost needs). For instance, Cursor lets you bring your own OpenAI key, or switch to Claude if you find it better for certain tasks. Supporting **open models** can also be a draw (e.g. Aider and Codex CLI letting you use local models – that pleases the open-source community).

Some startups choose to **open-source** their entire product (e.g. many CLI tools on GitHub, and even VS Code extensions like Cline). The reasoning: it's hard to compete with giants on data and model, but by being open you can become the community standard which the big players might contribute to (or at least not crush due to goodwill). We see this with the **MCP protocol** – started by open-source contributors, now adopted by Anthropic and Google, which is a big win for those standards. So a startup aligning with open standards can punch above its weight.

- **Developer Adoption Trends**: Developers have been quick to experiment with these tools, but adoption in production teams is still in early days. Copilot blazed the trail – by late 2023, over 50% of developers had tried some form of AI coding assistance [91] . Surveys show high satisfaction and perceived productivity boosts (e.g. Copilot users self-reported 55% faster coding for repetitive tasks) [92] . That opened managers' minds to integrating AI in the dev process. Now with agentic tools, we see a similar pattern: individual enthusiasts adopt them first, then team trials happen. Trust is a big factor – letting an AI *modify your codebase* or run code is a step further than just suggesting code. That's why tools like Codex CLI and Claude Code emphasize sandboxing and approval controls [93] [94] . Early team adopters often use these agents in "read-only" or suggestion modes until confidence grows. Over time, as success stories emerge (like "our AI agent caught a critical security bug in a PR" or "we built a microservice 2x faster with an AI pair programmer handling boilerplate"), comfort levels will rise.

Another trend: **multi-modal and multi-skill agents**. Developers appreciate when the AI can handle not just coding, but writing tests, generating docs, even creating diagrams or UI from descriptions. We see hints of this: Codex CLI and Gemini can take images (screenshots or design mockups) and generate code [95] [96] . Cascade can incorporate design assets with drag-and-drop images [52] . Claude Code and Gemini can browse the web for you to pull in library docs or StackOverflow answers on the fly [97] [98] . This Swiss-army-knife behavior is very appealing – it means the AI agent feels like an intern that can do many menial dev tasks (research, writing, testing) in one place. Developer communities (on Reddit, Discord) actively compare these tools' "skill sets" and share tips (for example, some swear Claude is better at understanding very large code contexts, whereas GPT-4 is better at precise logic; Cursor's small model is fast for autocomplete but they switch to GPT-4 for complex tasks). This cross-pollination is pushing all tools to broaden their capabilities.

- **Open vs Closed Ecosystems**: The battle is on between open ecosystems (where you can mix and match models and tools) and closed ones (end-to-end platforms). Big tech tends to push closed ecosystems: e.g. Microsoft wants you in VS Code + GitHub + Azure OpenAI; Google wants you using Code Assist + Cloud; Anthropic partnering with AWS for Bedrock integration. Startups and open-source devs push the mix-and-match approach: use an open CLI with whatever model API is cheapest that day, or plug one agent's extensions into another's interface. We see interesting collaborations, like *CodeRabbit integrating with Cursor and Windsurf* (they essentially said: use our AI reviews inside those editors, we don't mind if you use a different coding agent, we'll just focus on reviews [73] ). This suggests a potential **layered stack** emerging: you might use Windsurf IDE for coding with AI, CodeRabbit for reviewing AI, and your own scripts for deployment AI, etc. It's reminiscent of how DevOps tooling evolved – some prefer an all-in-one from a vendor, others compose best-of-breed open tools.

- **Monetization Strategies**: Many of these tools are currently free for individuals, which is amazing for devs right now. However, they will need to make money. Likely strategies:

    - **Subscription tiers** (e.g. Cursor Pro, Copilot paid, Windsurf Enterprise). The free tier hooks you; the paid tier gives more GPU-intensive features (like guaranteed GPT-4 access or higher rate limits).
    - **Enterprise licenses** – expect to see offerings like "run on-prem" or "private cloud deployment" for companies, at a premium. Codeium already had an on-prem option; CodeRabbit might offer a self-hosted reviewer for big companies concerned about code leaving their network.
    - **Model usage fees** – For open-source agents that let you bring your own model, they might partner with model providers. For instance, Google giving free Gemini now is surely to later convert some users to paid Vertex AI usage if they want scale or certain model versions.
    - **Marketplace or Platform model** – If an IDE like Windsurf or an agent like Claude Code becomes popular, they could have a plugin marketplace or paid add-ons. Imagine paying for a premium "security scan agent" extension or for priority support/feature requests.
    - **Data monetization** – This one is tricky given code privacy concerns, but anonymized analytics on how developers use these tools could inform model training (with user consent). Big providers likely use telemetry to improve the models (OpenAI said they might use Codex CLI interactions to refine their systems, similarly Google with Code Assist feedback).

Startups specifically need to justify valuations, so we might see them target management layers – e.g. an AI dashboard for team leads to see how much of the code was AI-generated and ensure quality (some companies, like Mendel.ai, are exploring AI governance tools). A startup could charge for those insights or compliance features built atop the agent usage.

## Future Trends: What's Next for AI Dev Stacks?

Looking ahead, several trends are poised to shape how we build software with AI:

- **AI Agent Teams**: We hinted at multi-agent earlier – this could become standard. Instead of one monolithic AI doing everything, you might have a **suite of specialized agents** working in concert. For example, an architecture agent designs a solution, a coding agent implements it, a testing agent generates tests, a devops agent sets up CI, and a project manager agent coordinates it all. It sounds futuristic, but all the pieces exist in isolation. The challenge is orchestration and trust between agents. Some research projects (and products like IBM's Watson Code Assistant in enterprise) are exploring defined roles for AI in the dev process. In practical terms, a near-future scenario: you file a ticket in Jira, an AI agent picks it up, writes code, another agent reviews it, and a human engineer just oversees and merges if satisfied. Humans may move more to **validation and creative design** roles, overseeing fleets of coding bots doing the grunt work.

- **Tighter IDE-Agent Fusion**: IDEs will likely morph to accommodate AI agents as first-class citizens. Expect more **"AI memory"** features in editors – e.g. the AI can store notes about why certain decisions were made, or pin important context (Windsurf's Memories is an early example [48] ). IDEs might offer *replay or diff views of AI actions*, so you can easily review what the agent did (some CLI tools already log all changes as git commits for traceability [97] ). Also, as trust grows, IDEs might

permit *"AI auto-pilot mode"* for certain safe tasks: e.g. one click and the AI will regenerate all outdated unit tests, commit them, and push.

- **Unified AI Dev Platforms**: It's possible we'll see consolidation where one platform handles coding, reviewing, and even **post-deployment monitoring with AI**. GitHub is well-positioned to try this (imagine Copilot not only helping you code, but also opening issues when production errors spike, and perhaps suggesting a fix commit). Similarly, cloud providers might integrate coding agents with cloud services: AWS's CodeWhisperer might tie into CodePipeline and CloudWatch (generate code + deploy infrastructure + watch for anomalies – all AI-assisted). This could streamline dev workflows but also raises the specter of **vendor lock-in** if all parts of your dev cycle rely on one provider's AI.

- **Open Source Counterweights**: The open-source community will continue producing alternatives to proprietary models, and could disrupt things if, say, someone releases a 70B ChatGPT-quality model that can run on a decent server. If devs can self-host an AI that's nearly as good as GPT-4 on code, many will do so for cost and privacy reasons. This might lead to more **community-driven AI agents** that rival the corporate ones. It's analogous to how Linux and open-source tools became standard in devops; we might see a world where an open "DevAgent" stack (open model + open agent software + custom tools) is the default for a lot of projects, with closed solutions mainly used by those who want turnkey convenience or specific premium features.

- **Monetization and Sustainability**: Eventually, the free lunches will end – either usage caps will tighten or features paywalled. There's also the question of **ethical and legal aspects**: AI-generated code can introduce licensing issues or security vulnerabilities. We may see services that **insure or certify AI-written code**, perhaps for a fee. Or regulatory compliance tools that verify AI code meets certain standards (imagine an AI that certifies "no GPL code was introduced" or "this code is not plagiarized from StackOverflow"). Startups might spawn around those assurances, especially for enterprise adoption where legal teams worry about AI contributions.

- **Developer Roles and Skills**: A softer trend is how this will change what being a developer means. If an AI agent can handle say 80% of the boilerplate and rote coding, developers might spend more time on creative problem decomposition, validating AI output, and higher-level thinking. **"Prompt engineering"** becomes a bit of a developer skill – knowing how to instruct AI effectively. But unlike generic prompt engineering, in coding the prompts are often just normal language descriptions of features or bugs, so it's more about **spec writing and critical reading** of AI output. Early adopters report that using these agents forces them to **articulate requirements more clearly** and think in terms of step-by-step plans to guide the AI. In a way, it's making explicit the thought process that was implicit. Junior devs might get up to speed faster with an AI pair programmer always there to explain code or suggest approaches. Senior devs might leverage agents to handle tedious refactors, freeing them to focus on architecture. Roles like QA and DevOps might shift too – if AI can generate tests and manage infra configs, those specialists might move into supervising those AI outputs and focusing on edge cases the AI doesn't handle. All told, the **human element remains crucial** – but the day-to-day tasks will likely change. It's reminiscent of the shift when high-level languages or frameworks came out: less fiddling with low-level code, more focus on design and integration.

# Key Takeaways for Developers, Teams, and Builders

For individual devs and small teams, this landscape can be both exciting and overwhelming. Here are some takeaways and tips on navigating the new AI dev stack:

- **Leverage the Free Tools to Boost Productivity:** If you haven't already, try out at least one agentic AI coding tool in your workflow. For instance, use GitHub Copilot (or Cursor's free tier) for code suggestions and documentation, and try a terminal agent like OpenAI's Codex CLI or Claude Code for automating those "glue" tasks (running tests, applying large-scale refactors, etc.). These tools can significantly speed up routine tasks – users report double-digit percent gains in productivity and far less frustration on boilerplate code 92 72 . Since many offerings are currently free or have generous trials (Gemini CLI's effectively unlimited personal use, Windsurf free editor, CodeRabbit free in editor, etc.), you can experiment at low cost. Just remember to secure any API keys and be mindful of sending proprietary code to cloud models if you're not comfortable with that.

- **Pick the Right Tool for the Job:** The landscape isn't either/or – you can mix tools. **Terminal agents** like Claude Code or Gemini CLI are fantastic for *project-wide operations* (e.g. "upgrade all our dependencies and fix breaking changes" or "find and fix all TODO comments in the repo"). They shine when you need to affect many files or use external tools autonomously. **IDE agents** like Cursor or Copilot are great for *in-the-flow coding* – getting suggestions as you type, quick fixes, or asking "what does this function do?". And **code review bots** like CodeRabbit or Ellipsis help maintain quality after the code is written. An early-stage startup, for example, might use Copilot in VS Code to scaffold features, Claude Code in the terminal to handle repo-wide refactors or generating config files, and CodeRabbit to double-check PRs for mistakes. It's not as heavy as it sounds – each integrates in context (Copilot and CodeRabbit in GitHub, Claude in CLI) without a lot of overhead. Identify your pain points (Is writing tests slow? Is conforming to style tedious? Are you unsure about architecture?) and bring in the AI tool that best addresses each.

- **Maintain Human Oversight and Best Practices:** These agents are powerful, but they are not infallible. They can and **will** make mistakes – an AI might introduce a subtle bug, mis-interpret a requirement, or produce code that just doesn't idiomatically fit your style. Always **review AI-generated code** (just as you would a junior developer's code). Use your version control to see diffs of what the agent changed. Write additional tests if needed to validate its outputs. Think of the AI as an apprentice – fast and tireless, but needing guidance. Also, continue to do code reviews among humans even if an AI did one – you'll catch things the AI missed and vice versa. Importantly, integrate these tools into your **existing processes** rather than letting them dictate process. For example, if you require design docs before implementation, use the AI to help write the doc, but don't skip the design approval just because the AI can start coding immediately. *You still set the standards*; the AI just helps you meet them more efficiently.

- **Be Mindful of Data and Privacy:** When using cloud-based AI (OpenAI, Anthropic, etc.), remember that your code is being sent to their servers. Most companies assure they don't train on your data (OpenAI has a policy for API data, CodeWhisperer doesn't retain code, etc.), but it's still a third party handling your IP. If your project is super sensitive or proprietary, consider using tools that allow local models or self-hosting. For instance, run Codex CLI with a local LLM for those parts, or use an open-source model behind Aider. The performance might be lower, so you might use a hybrid approach (do initial development with cloud AI, but maybe don't feed it your entire proprietary codebase

history). On the flip side, if you *are* comfortable (or working on open source code anyway), take advantage of the superior capabilities of models like GPT-4 – just ensure any credentials or personal data are excluded from prompts (agents like Claude Code have settings to mask things like secrets, but double-check!).

- **Customize and Extend if You Can:** Many of these tools allow **customization** that can give you an edge. For example, try writing a **GEMINI.md or CLAUDE.md** file at your repo root with notes for the agent (model-specific config): you can put high-level instructions like "This is a fintech app, user data must be validated with these rules, prefer functional style code" – the agents will read that and follow consistently (both Claude and Gemini support such project-wide context files via MCP and system prompts) [99] . If you have internal tools or APIs, consider hooking them in via an MCP server or custom function – e.g. you could let the agent call a `getDesignDocs()` function to retrieve internal documentation when needed. This does require some coding, but the SDKs are there [83] [84] and even a basic integration can save time (imagine the agent automatically pulling the API contract for a microservice you're calling, instead of hallucinating the interface). By tuning the AI's environment to your project, you turn it from a generic assistant into a bespoke team member.

- **Educate and Set Team Policies:** For teams, it's worth having a discussion about **how to use AI assistants responsibly**. Decide on things like: Should AI-generated code be marked or commented somehow (at least in commit messages, e.g. "Co-authored-by: AI")? Are there parts of the codebase or tasks off-limits to AI (maybe critical security code that only certain senior devs handle)? How will code reviews account for AI contributions – do you require an extra manual review if an agent was used heavily on a commit? Also share best practices among the team – one dev might find that "Claude Code is great for research, but we only trust its code after thorough testing" or that "Copilot's suggestions for our Python code need careful type checking." Spreading this knowledge will help everyone use the tools effectively and avoid pitfalls. **Mentorship can incorporate AI** too: senior devs guiding juniors not only in coding but in how to get the best out of AI tools, how to verify AI output, etc. The tools are new, so fostering a culture of learning and safety around them is important.

- **Stay Adaptable:** The AI tooling landscape is changing quickly – what's state-of-art now might be eclipsed in six months. Remain adaptable in your stack. Avoid overly binding your workflow to one proprietary solution if possible. For instance, if you adopt an AI agent, ensure you can export or retain the code it produces independent of the tool (which you generally can, since it's just your code repository). Keep an eye on emerging tools – perhaps set aside time every quarter to re-evaluate: "Is there a new open-source model that we can run cheaper?" or "Does our current tool still give us a competitive edge or has it become table stakes?" Many teams are even creating an **"AI tools" team or champion** – someone who evaluates new releases (like GPT-5 when it comes, or Gemini updates, etc.) and figures out if they can benefit the team. Given that productivity gains from AI could become a differentiator in software delivery, it's wise to stay on the curve. At the same time, don't chase every shiny object – establish some criteria (does it support our programming languages? our security requirements? etc.) and be strategic.

- **Embrace the Change**: Finally, attitude matters. There is understandable anxiety about AI in programming – fears of job impact or deskilling. But history with developer tools (from compilers to StackOverflow) shows that those who *embrace automation and abstraction* tend to accelerate their careers, not lose them. Let the repetitive parts of coding be handled by AI – this frees you to focus on

creative engineering: designing better systems, solving novel problems, and orchestrating technology in ways that add real value. Use the AI to amplify your strengths and to shore up your weaker areas (not great at front-end UI? an AI can draft the HTML/CSS; hate writing tests? let the AI stub them and you refine). By integrating these agents into your workflow, you essentially become a **tech lead managing a team of AI juniors**. Your role elevates to supervisor and architect, which is a valuable position. And for those just learning to code, don't worry – using AI doesn't mean you won't learn; in fact, with proper use, it can accelerate learning (the AI explains things, suggests idiomatic patterns, etc.). Just be sure to **actively engage** rather than blindly accept outputs, and you'll find you pick up skills faster with an ever-present tutor/assistant.

---

**In conclusion**, the landscape of AI coding tools in late 2025 is vibrant and quickly maturing. We have powerful terminal agents that can take a natural language goal and turn it into committed code [2] [10] . We have smart IDE assistants that deeply understand our projects and collaborate in real-time [35] [44] . We have bots that guard our code quality in reviews and automate away nitpicks [65] [66] . And we have emerging frameworks to extend these capabilities even further, potentially knitting multiple agents together. The savvy developer or team will take advantage of these tools to **build faster and smarter**, while also putting in place the right practices to maintain quality and control.

The **"agentic stack"** – your chosen combo of AI coding agent, code reviewer, and custom tooling – could become as important as your programming language or cloud provider in determining productivity. It's an exciting time to be building software, as we stand on the cusp of a new paradigm where human creativity pairs with artificial intelligence to produce software in a more interactive, high-level way. Those who experiment and adapt will find that far from replacing developers, these AI agents are *augmenting developers* – letting us focus on the fun parts of coding while the grunt work is handled by tireless silicon colleagues. In the end, coding has always been about using the best tools to get from idea to reality; agentic AI is simply the next leap in our toolset. Happy coding (with your new AI friends)!

**Sources:**

- Anthropic Claude Code Overview [10] [11]
- Google's Gemini CLI Announcement [21] [14]
- Windsurf (Codeium) Cascade Features [44] [51]
- Cursor AI Editor – Daily.dev report [37] [35]
- OpenAI Codex CLI – GitHub README [2] [3]
- CodeRabbit – Google Cloud Blog [100] [66]
- Developer Adoption Stats – Reddit/GitHub Blog [1] [92]

---

[1] Coding is no more a MOAT. 46% of codes on GitHub is already being built using GitHub Copilot across all programming languages : r/singularity

https://www.reddit.com/r/singularity/comments/1177kwz/coding_is_no_more_a_moat_46_of_codes_on_github_is/

[2] [7] [8] [86] [87] [94] GitHub - openai/codex: Lightweight coding agent that runs in your terminal

https://github.com/openai/codex

[3] [4] [5] [6] [93] [95] OpenAI Codex CLI: Build Faster Code Right From Your Terminal | Blott Studio

https://www.blott.studio/blog/post/openai-codex-cli-build-faster-code-right-from-your-terminal

[9] [10] [11] [12] [13] [97] Claude Code overview - Anthropic

https://docs.anthropic.com/en/docs/claude-code/overview

[14] [16] [17] [18] [19] [20] [21] [22] [23] [24] [25] [26] [59] [60] [61] [98] [99] Google announces Gemini CLI: your open-source AI agent

https://blog.google/technology/developers/introducing-gemini-cli-open-source-ai-agent/

[15] Model Context Protocol (MCP) - Anthropic

https://docs.anthropic.com/en/docs/claude-code/mcp

[27] aider is AI pair programming in your terminal - GitHub

https://github.com/lloydchang/Aider-AI-aider-fka-paul-gauthier-aider

[28] What is Aider? : r/LocalLLaMA - Reddit

https://www.reddit.com/r/LocalLLaMA/comments/1izm42j/what_is_aider/

[29] Aider: AI pair programming in your terminal | Hacker News

https://news.ycombinator.com/item?id=39995725

[30] Cline - AI Autonomous Coding Agent for VS Code

https://cline.bot/

[31] [32] Devin AI review: is it better than Cursor?

https://www.builder.io/blog/devin-vs-cursor

[33] [34] [35] [36] [37] [38] [39] [40] [41] Cursor AI: The AI-powered code editor changing the game

https://daily.dev/blog/cursor-ai-everything-you-should-know-about-the-new-ai-code-editor-in-one-place

[42] Windsurf Editor | Windsurf (formerly Codeium)

https://windsurf.com/editor

[43] [46] [48] [51] [52] [53] [55] [56] [57] [58] Windsurf (formerly Codeium) - The most powerful AI Code Editor

https://windsurf.com/

[44] [45] [47] [49] [50] [96] Windsurf AI Agentic Code Editor: Features, Setup, and Use Cases | DataCamp

https://www.datacamp.com/tutorial/windsurf-ai-agentic-code-editor

[54] Windsurf just destroyed AI coding agents with something far far better | by Tari Ibaba | Coding Beauty | May, 2025 | Medium

https://medium.com/coding-beauty/new-windsurf-swe-agents-e9658eab0939

[62] [63] [64] [65] [66] [67] [68] [71] [76] [77] [100] How CodeRabbit built its AI code review agent with Google Cloud Run | Google Cloud Blog

https://cloud.google.com/blog/products/ai-machine-learning/how-coderabbit-built-its-ai-code-review-agent-with-google-cloud-run

[69] [70] [72] [74] [75] [81] [82] AI Code Reviews | CodeRabbit | Try for Free

https://www.coderabbit.ai/

[73] CodeRabbit's AI Code Reviews Now Live Free in VS Code, Cursor

https://thenewstack.io/coderabbits-ai-code-reviews-now-live-free-in-vs-code-cursor/

[78] AI Code Reviews | CodeRabbit | Try for Free

https://coderabbit.ai/

[79] Code Review - Ellipsis | AI Code Reviews & Bug Fixes - Introduction

https://docs.ellipsis.dev/features/code-review

[80] Ellipsis: AI code reviews & bug fixes | Y Combinator

https://www.ycombinator.com/companies/ellipsis

[83] [84] [85] [90] Claude Code SDK - Anthropic

https://docs.anthropic.com/en/docs/claude-code/sdk

[88] [89] Devon | AI Agents Directory

https://aiagentslist.com/agent/devon

[91] Coding on Copilot: 2023 Data Suggests Downward Pressure on …

https://www.gitclear.com/coding_on_copilot_data_shows_ais_downward_pressure_on_code_quality

[92] GitHub Copilot · Your AI pair programmer

https://github.com/features/copilot