## Reviewing the Landscape of AI Tooling: Hype Justified

The integration of artificial intelligence into various technological tools has moved beyond theoretical discussions, establishing itself as a tangible force reshaping how individuals and organizations approach their work. The initial assessment that AI tooling is "worth the hype" resonates with the observed advancements and the potential these tools hold for enhancing productivity and fostering innovation. This report delves into a selection of AI-powered solutions, spanning search engines, locally hosted models, and integrated development environment (IDE) plugins, to examine the validity of this enthusiastic outlook.

## Reimagining Search: The Dawn of Al-Powered Discovery

The way individuals seek information is undergoing a significant transformation, moving from traditional keyword-based searches to more intuitive, AI-driven discovery. This shift is characterized by platforms that aim to understand the intent behind user queries and provide direct, synthesized answers accompanied by verifiable sources. The emergence of multiple AI search engines signifies a growing recognition of this evolving need and a competitive landscape striving to deliver more intelligent and efficient information retrieval.

## Perplexity AI: Beyond Traditional Search

Launched in 2022, Perplexity AI has rapidly gained popularity among a diverse user base, including researchers and professionals, due to its capacity to comprehend context, summarize information, and present reliable sources in real-time.<sup>1</sup> Unlike conventional search engines that primarily return a list of links, Perplexity leverages sophisticated Natural Language Processing (NLP) and machine learning to interpret queries with human-like understanding.<sup>1</sup> This allows the platform to break down complex questions and generate clear, informative responses.<sup>1</sup> A key differentiator of Perplexity is its commitment to transparency, as every answer includes links to the original content it draws from, enabling users to independently verify the results.<sup>1</sup> This feature is particularly valuable in academic research, where the ability to trace facts back to reputable sources like academic journals and news outlets is crucial.<sup>1</sup>

Perplexity offers a range of features designed to streamline the information-gathering process. It can summarize information from multiple sources into concise overviews, saving users the effort of scanning numerous pages.<sup>1</sup> The platform also supports contextual dialogue, allowing users to ask clarifying or related follow-up questions without having to start their search anew.<sup>1</sup> Answers are generated from live web data, ensuring that users receive up-to-date information rather than relying solely on pre-trained models.<sup>1</sup> While basic usage does not require an account, creating one unlocks additional benefits such as chat history,

saved threads, and access to Pro features, including advanced AI models and priority responses.<sup>1</sup>

The applications of Perplexity AI extend across various domains. In academic research, it proves useful for literature reviews, sourcing recent studies, summarizing dense academic concepts, and quickly understanding opposing viewpoints.<sup>1</sup> Professionals across industries utilize it for tasks like conducting market research, exploring industry trends, answering technical or legal questions, and summarizing lengthy reports or articles.<sup>1</sup> The platform offers different search modes to cater to varying needs. "Quick Search" provides immediate responses by summarizing information from indexed sources, while "Pro Search" delves deeper, potentially asking follow-up questions to refine the search and deliver more tailored answers.<sup>3</sup> For even more complex inquiries, the "Deep Research" feature conducts extensive searches, reads numerous sources, and reasons through the material to autonomously generate comprehensive reports, suitable for expert-level analysis in fields like finance, marketing, and technology.<sup>4</sup> This capability significantly reduces the time required for in-depth research, performing tasks that would typically take a human expert many hours in just a few minutes.<sup>4</sup>

Comparisons with traditional search engines like Google highlight Perplexity's strengths in providing direct answers with clear citations, offering a cleaner, ad-free interface.<sup>5</sup> Unlike Google, which often presents a list of ranked websites, Perplexity prioritizes concise, sourced responses, making it particularly valuable for research where accuracy and depth are paramount.<sup>7</sup> While Perplexity generally strives for accuracy, it has been noted that occasionally, generated answers might include "extra details" that are not entirely factual <sup>2</sup>, underscoring the importance of critical evaluation. Nevertheless, its focus on providing transparent, source-backed answers and its ability to handle complex queries make Perplexity AI a powerful tool for efficient and reliable information discovery.

## Phind: The Developer's Intelligent Search Companion

Phind emerges as an AI-powered search engine specifically tailored to meet the information needs of developers and other professionals seeking accurate and efficient solutions to coding challenges and technical inquiries.<sup>12</sup> Unlike general-purpose search engines, Phind leverages advanced AI technologies to understand and respond to complex technical queries, providing precise and contextually relevant answers.<sup>12</sup> This focus on developer-centric needs is evident in its key features, which are designed to streamline the search process and integrate seamlessly into the developer workflow.

A significant aspect of Phind is its integration with popular development environments, such as Visual Studio Code (VS Code).<sup>12</sup> This integration allows developers to access information directly within their coding environment, facilitating quick access to coding solutions, technical documentation, and useful code snippets without the disruption of switching contexts.<sup>12</sup> Phind employs sophisticated AI algorithms to deliver precise, context-aware search results tailored for developers.<sup>12</sup> It processes complex queries in real-time, ensuring users receive accurate and contextually appropriate answers swiftly.<sup>12</sup> The platform also

provides instant access to a vast array of technical documentation and coding solutions, aiding in rapid problem-solving.<sup>12</sup> Furthermore, Phind offers collaboration tools that enable users to share and discuss search results, enhancing teamwork and knowledge sharing among team members.<sup>12</sup>

Phind stands out for its fast performance, utilizing advanced models like Phind-70B, which is based on the CodeLlama-70B architecture and fine-tuned on an extensive dataset.<sup>13</sup> This model supports a large context window, allowing for more comprehensive analysis and responses to complex programming tasks.<sup>13</sup> The platform can generate, test, and run code directly in the browser, further enhancing the coding workflow and productivity.<sup>13</sup> Phind supports a wide range of programming languages, frameworks, databases, and cloud platforms, making it a versatile tool for developers working across different technological ecosystems.<sup>14</sup> Notably, it can integrate with users' existing codebases to provide tailored debugging and development recommendations using advanced AI technologies.<sup>14</sup> Unlike traditional search engines that primarily return links, Phind focuses on understanding user intent and providing actionable answers.<sup>16</sup> It can even present answers visually with inline images, diagrams, and interactive widgets to enhance understanding.<sup>17</sup> While both Perplexity and Phind are AI-enhanced search engines, Phind's stronger emphasis on technical and coding-related information distinguishes it as a specialized tool for developers.<sup>18</sup> Its ability to process large amounts of information, generate clean and optimized code, and assist with debugging tasks makes it a valuable asset in the software development process.<sup>16</sup>

#### Gemini: Google's Al Ecosystem and Search Enhancement

Gemini represents Google's most advanced AI model, which is being integrated across its suite of products and services to enhance user experience and capabilities.<sup>19</sup> Evolving beyond a standalone model, Gemini is forming a powerful ecosystem that aims to make Google's offerings faster, smarter, and more helpful.<sup>19</sup> This integration spans a wide array of applications, from search and email to document creation and collaboration tools. Within Google Search, Gemini powers "Al Overviews," which provide quick, Al-generated summaries at the top of search results, offering a conversational way to get answers and ask follow-up questions.<sup>19</sup> This feature aims to streamline information retrieval by presenting key insights directly, saving users the need to sift through multiple links.<sup>21</sup> Gemini's capabilities extend to other Google Workspace applications, with features like "Help me write" in Gmail and Docs assisting users in drafting content, correcting grammar, and translating text.<sup>19</sup> It can also aid in creative tasks, such as generating images and designs in Slides.<sup>19</sup> Gemini offers advanced functionalities for more complex tasks. "Deep Research" allows users to research and synthesize information from the web, saving time by acting as a personal AI research assistant.<sup>23</sup> Gemini can now connect with users' Google apps and services to provide more personalized responses, referencing recent searches or travel plans, for example.<sup>24</sup> Users can even create custom "Gems" to personalize Gemini for specific tasks, such as translation or meal planning.<sup>24</sup> Built from the ground up for multimodality, Gemini can process and reason across various formats, including text, code, images, and video.<sup>19</sup> Different

versions of the Gemini model, such as Ultra, Pro, and Nano, are optimized for different purposes, ranging from highly complex tasks to efficient on-device processing.<sup>27</sup> Gemini's integration with Google Assistant enables voice commands and context awareness based on what's on the user's screen on Android devices.<sup>26</sup> This allows for seamless interaction and assistance across different contexts. While Gemini aims to enhance the search experience, some users still find the traditional Google Search more suitable for certain types of queries, particularly those requiring quick lookups of local information.<sup>29</sup> The underlying technology of Gemini involves training on massive amounts of data, enabling it to understand and generate responses, which differs fundamentally from Google's traditional search algorithms that index and rank existing web content.<sup>27</sup> Despite some user preferences for the conventional search interface, Gemini represents a significant step in Google's strategy to embed advanced AI capabilities deeply within its ecosystem, offering a versatile AI assistant that extends far beyond traditional search functionalities.

## Unleashing Local Al Power: Privacy and Control in Your Hands

A growing trend in the AI landscape is the development of tools that allow users to run AI models locally on their own hardware. This approach offers several compelling benefits, including enhanced privacy by keeping data on-device, greater control over the AI's operation, and the potential for offline access, independent of internet connectivity. The emergence of various local AI tools indicates a strong interest in harnessing AI capabilities directly on personal computers and within private networks.

## Pieces: Streamlining Development Workflows with CLI Intelligence

Pieces presents itself as a comprehensive tool designed to enhance developer workflows through intelligent command-line interface (CLI) interactions with its underlying Pieces OS.<sup>33</sup> This CLI agent provides a wide array of functionalities aimed at managing development-related assets and integrating AI capabilities directly into the terminal. The Pieces CLI agent offers robust asset management features, allowing users to list, modify, create, delete, and open various digital materials within their Pieces Drive.<sup>33</sup> It also facilitates interaction with registered applications and provides control over AI models, enabling users to list available models and select the one to be used for specific tasks.<sup>33</sup> Developers can even execute Pieces bash materials directly through the CLI.<sup>33</sup> For configuration, the tool allows users to view the current settings and set their preferred editor.<sup>33</sup> A core functionality is the ability to ask questions to AI models, with the option to include specific materials or files/folders as context to guide the AI's responses.<sup>33</sup>

The CLI also includes comprehensive chat management features, allowing users to list existing chats, view messages within a chat, switch between different conversations, create new chats, delete the current chat, and rename ongoing discussions.<sup>33</sup> For version control, Pieces integrates with GitHub, enabling users to commit changes and automatically generate

commit messages, with the option to push these changes directly.<sup>33</sup> The tool also provides powerful search capabilities, including fuzzy search for general queries, neural code search for code-specific inquiries, and full-text search for broader content discovery.<sup>33</sup> Basic functionalities like login and logout are also supported, along with commands to retrieve version information, access help documentation, initiate the onboarding process, send feedback, and find contribution guidelines.<sup>33</sup> Additionally, the Pieces CLI can even trigger the installation of the Pieces OS and clear the terminal interface.<sup>33</sup>

The potential applications of the Pieces CLI agent are diverse, catering to developers who prefer the efficiency and automation of the command line for managing code snippets, notes, and other development assets.<sup>33</sup> Content creators can leverage it to organize and access various forms of content, while researchers can manage research notes, papers, and data.<sup>33</sup> The ability to interact with AI models for coding assistance, content generation, information retrieval, or analysis further enhances its utility across these domains.<sup>33</sup> In essence, Pieces provides a powerful and efficient way to interact with its ecosystem without relying on a graphical user interface, appealing to users who value the speed and automation capabilities of the command line.

#### Goose: Automating Complex Development Tasks Autonomously

Goose represents a significant advancement in AI tooling, offering an open-source AI agent designed to autonomously handle complex development tasks from initiation to completion.<sup>34</sup> Going beyond the capabilities of simple code suggestion tools, Goose can build entire projects from scratch, write and execute code, debug failures, orchestrate intricate workflows, and interact with external APIs without direct human intervention.<sup>34</sup>

The primary goal of Goose is to empower developers to accelerate their work and concentrate on innovation by automating time-consuming and repetitive aspects of the development process.<sup>34</sup> Whether it's rapidly prototyping a new idea, refining existing codebases, or managing complex engineering pipelines, Goose is designed to adapt to the user's workflow and execute tasks with precision.<sup>34</sup> Its architecture prioritizes maximum flexibility, allowing it to work with any Large Language Model (LLM) and seamlessly integrate with MCP servers.<sup>34</sup> Goose is available in two forms: a desktop application for users who prefer a graphical interface and a CLI for those who favor command-line interactions.<sup>34</sup>

The potential applications of Goose are vast for developers seeking to enhance their productivity and focus on higher-level strategic work.<sup>34</sup> By taking over the more mechanical aspects of development, such as generating boilerplate code, debugging common errors, and managing project dependencies, Goose can free up developers to dedicate their time and energy to more creative problem-solving and the exploration of new technologies.<sup>34</sup> Its ability to interact with external APIs autonomously opens up possibilities for integrating various services and automating complex integrations, further streamlining development workflows.<sup>34</sup> Goose signifies a move towards a future where AI agents can play a more active and autonomous role in the software development lifecycle, potentially leading to significant gains in efficiency and speed.

## aichat: A Universal Interface for Diverse Language Models

aichat distinguishes itself as a versatile, all-in-one Large Language Model (LLM) CLI tool that provides a unified interface for interacting with a vast array of language models from over 20 leading providers.<sup>35</sup> This extensive multi-provider support allows users to leverage the strengths of different LLMs through a single, consistent command-line experience. The core of aichat lies in its powerful command-line interface (CMD mode) and its interactive Chat-REPL (Read-Eval-Print Loop) mode.<sup>35</sup> The REPL mode enhances the conversational experience with features like tab autocompletion, support for multi-line input, history search, configurable keybindings, and customizable REPL prompts.<sup>35</sup> For users who frequently work with the shell, aichat offers a Shell Assistant that can translate natural language descriptions of tasks into precise shell commands, adapting to the user's operating system and shell environment to improve command-line efficiency.<sup>35</sup> The tool is designed to handle diverse input formats, including stdin, local files and directories, and remote URLs, providing flexibility in how data is processed.<sup>35</sup> It can also incorporate the last reply, outputs from external commands, and combine various input types for more complex interactions.<sup>35</sup> aichat allows for extensive customization through role definitions, where users can tailor the LLM's behavior for specific interactions by defining prompts and model configurations.<sup>35</sup> It also supports context-aware conversations through session management, ensuring continuity across multiple turns.<sup>35</sup> For repetitive tasks, users can define macros by combining a series of REPL commands into custom shortcuts.<sup>35</sup> A particularly powerful feature is Retrieval-Augmented Generation (RAG), which enables aichat to integrate external documents into LLM conversations, grounding the model in specific knowledge for more accurate and contextually relevant responses.<sup>35</sup> Furthermore, aichat supports function calling, allowing LLMs to connect to external tools and data sources, expanding their capabilities beyond core functionalities to tackle a wider range of tasks. This includes the ability to integrate AI tools for automation and the creation of AI Agents defined by instructions, tools, and documents.<sup>35</sup> Beyond its CLI capabilities, aichat includes a built-in lightweight HTTP server for easy deployment, offering APIs for Chat Completions, Embeddings, and Rerank, as well as a web-based LLM Playground and LLM Arena for comparing different LLMs.<sup>35</sup> The tool also supports custom dark and light themes to enhance the visual experience.<sup>35</sup> Due to its diverse features, aichat can be used in numerous ways, from general chat and command-line task automation to data analysis, code generation, information retrieval, and building sophisticated Al agents. Its ability to act as a universal interface for a wide range of LLM providers makes it a valuable tool for users who want to explore and leverage the diverse capabilities of the current Al landscape.

#### Jan: Your Private Al Playground, Offline and Secure

Jan positions itself as a user-friendly, open-source alternative to ChatGPT, with a key focus on operating entirely offline on the user's computer.<sup>36</sup> This offline functionality provides users with complete control over their AI interactions and ensures the privacy of their data.

A central feature of Jan is its integrated Model Library, which includes a variety of popular Large Language Models (LLMs) such as Llama, Gemma, Mistral, and Qwen.<sup>36</sup> This allows users to easily access and run different models based on their specific needs and preferences. While primarily designed for offline operation, Jan also offers the flexibility to connect to remote AI APIs like Groq and OpenRouter, providing users with access to cloud-based models if desired.<sup>36</sup> Furthermore, Jan includes a local API server that is designed to be equivalent to the OpenAI API. This compatibility enables seamless integration with other applications that are built to interact with OpenAI's models.<sup>36</sup>

Jan supports the use of Extensions, allowing users to customize and enhance its functionalities beyond the core features.<sup>36</sup> The platform and its underlying engine, Cortex, are designed for broad compatibility, supporting various architectures, including NVIDIA GPUs, Apple M-series and Intel processors, as well as Linux and Windows operating systems.<sup>36</sup> The target user for Jan is described as a "layperson" who seeks a straightforward way to download and run LLMs and utilize AI with complete control and privacy.<sup>36</sup> This emphasis on user-friendliness and privacy suggests that Jan is designed to be accessible to individuals who may not have extensive technical expertise in AI. The focus on offline operation makes it particularly appealing to users who are concerned about data security or who need to use AI in environments with limited or no internet access. By providing a private and locally hosted AI experience, Jan democratizes access to powerful language models and empowers users to explore the potential of AI on their own terms.

## Elevating the Coding Experience: Al Assistants in Integrated Development Environments

The integration of AI directly into Integrated Development Environments (IDEs) represents a significant step towards enhancing coding productivity and streamlining the software development process. These AI assistants, often in the form of plugins, offer a range of features designed to help developers write code faster, identify errors more easily, and understand complex codebases more effectively. The increasing number of such plugins highlights the growing recognition of the value that AI can bring to the daily tasks of software developers.

## ProxyAI: Intelligent Code Assistance for JetBrains

ProxyAI emerges as an open-source AI copilot specifically designed for the JetBrains suite of IDEs, aiming to provide a robust alternative to other AI-powered code assistants available in the market.<sup>37</sup> This plugin offers a comprehensive set of features that span both chat-based interactions and direct code manipulation, all intended to enhance the developer experience within the familiar JetBrains environment.

ProxyAI's chat features are designed to provide developers with contextual assistance and quick access to relevant information. The "Auto Apply" feature allows developers to directly integrate AI-suggested code changes into their editor, providing a diff view for preview and a one-click acceptance or rejection mechanism.<sup>37</sup> Users can engage in chat with the AI using

images as context, either by manually uploading them or allowing ProxyAI to automatically detect screenshots, adding a visual dimension to the assistance.<sup>37</sup> The plugin also enables developers to quickly access and reference their project files and folders within the chat session, allowing the AI to provide more context-aware coding suggestions.<sup>37</sup> Furthermore, ProxyAI facilitates referencing external web documentation, such as API guides, and the project's Git history directly within the chat, aiding in quick answers and understanding code evolution.<sup>37</sup> The AI can also connect to the internet through the chosen Large Language Model (LLM), allowing it to search for the most up-to-date information to answer developer questions.<sup>37</sup> For a more tailored experience, users can choose from multiple different personas for the AI assistant based on their specific needs, whether it's learning, code writing, or proofreading.<sup>37</sup>

In terms of code-specific features, ProxyAl offers "Next edits," which provides multi-line code edit suggestions based on recent activity as the developer types, anticipating their needs and speeding up the process.<sup>37</sup> It also offers both single-line and whole-function autocomplete suggestions to help write code faster and with fewer errors.<sup>37</sup> A particularly intuitive feature is the ability to edit code using natural language: developers can highlight the code they want to modify and describe the desired changes in plain English, with ProxyAl attempting to apply these modifications.<sup>37</sup> The plugin can also provide context-aware naming suggestions for methods, variables, and other code elements, promoting clarity and consistency.<sup>37</sup> Finally, ProxyAl can generate concise and descriptive commit messages based on the changes made in the codebase, streamlining the version control process.<sup>37</sup> These features collectively contribute to significant benefits for developers, including enhanced productivity, improved code quality, faster problem-solving, streamlined workflows, greater flexibility due to its open-source nature and model agnosticism, and personalized assistance through customizable personas.

## Tabby: The Open-Source, Self-Hosted AI Coding Companion

Tabby positions itself as an open-source, self-hosted AI coding assistant that serves as a compelling alternative to GitHub Copilot, offering developers greater control and transparency over their AI-powered coding assistance.<sup>38</sup> Its self-contained nature means it does not require a database management system (DBMS) or cloud services to operate.<sup>39</sup> Tabby features an OpenAPI interface, making it easy to integrate with existing infrastructure, such as cloud-based IDEs.<sup>39</sup> Notably, it supports consumer-grade GPUs, allowing for efficient operation on readily available hardware.<sup>39</sup>

A core feature of Tabby is its advanced code understanding, which enables it to offer precise autocompletion suggestions by analyzing the context within a codebase, thereby saving developers time and effort and improving code quality.<sup>38</sup> Being fully open source, Tabby allows developers not only to use the tool but also to modify it to better suit their specific needs, whether it's adding new features, tweaking the underlying model, or integrating it with other tools.<sup>38</sup> Tabby is designed for seamless integration with popular IDEs and text editors, including Visual Studio Code, JetBrains, and Vim/NeoVim, ensuring that developers can

continue working in their preferred environment without disruption.<sup>38</sup> It provides real-time code suggestions as the developer types, helping to write code faster and more accurately.<sup>38</sup> Tabby's design philosophy emphasizes an end-to-end approach, optimizing not just the interaction with Large Language Models but also the IDE extensions to ensure rapid and accurate code completion.<sup>40</sup> It achieves accurate streaming and cancellation with an adaptive caching strategy to ensure completion times of less than a second.<sup>40</sup> Tabby also parses relevant code into Tree Sitter tags to provide more effective prompts to the underlying language models.<sup>40</sup> While primarily focused on code completion, some sources suggest that a product also named "Tabby" by BytePlus offers features like intelligent code generation, error detection, integration with business analytics, and scalable API access.<sup>42</sup> It's important to note that these might refer to different products with the same name. Regardless, the open-source Tabby stands out as a secure, flexible, and transparent AI coding assistant that developers can self-host and customize, providing a compelling alternative to proprietary solutions.

### Privy: Local Al Autocomplete and Chat within VSCode

Privy is a Visual Studio Code (VS Code) extension that brings the power of AI-driven code autocomplete and chat directly to the developer's local machine.<sup>44</sup> A key aspect of Privy is its ability to run Large Language Models (LLMs) locally using platforms like Ollama, llamafile, or llama.cpp, ensuring that code suggestions and chat interactions are processed privately and without the need for cloud connectivity.<sup>44</sup>

Privy offers a range of features designed to enhance the coding experience within VS Code. It provides intelligent code suggestions as the developer types, helping to speed up the coding process and reduce errors.<sup>44</sup> The extension also includes a chat interface that allows developers to interact with the locally running LLM in a copilot-style manner, enabling them to ask questions about their code, explain complex logic, generate unit tests, find potential bugs, and diagnose errors.<sup>44</sup> Conversations within the chat are threaded, helping to keep different topics organized and easy to follow.<sup>44</sup> Being open source, Privy emphasizes privacy, ensuring that sensitive code and conversations remain on the user's local system.<sup>44</sup> It supports various popular LLMs, including DeepSeek Coder, CodeLLama, and Mistral, giving developers the flexibility to choose the model that best suits their needs and hardware capabilities.<sup>44</sup> To get the most out of Privy, it's recommended that users be specific and provide sufficient context when asking for assistance.<sup>44</sup> While Privy is a significant step towards having a local AI coding assistant, it's important to note that, like any AI tool, it may occasionally provide inaccurate answers, especially on less common topics or during detailed conversations.<sup>44</sup> Therefore, developers should exercise caution and not blindly trust the AI's responses without verification.<sup>44</sup> Privy aims to provide a secure and private AI-powered coding assistant within VS Code, leveraging the power of local LLMs to enhance productivity and offer valuable coding support.

## Ollama Autocoder: Seamless Local Code Completion with Ollama

Ollama Autocoder is a Visual Studio Code (VS Code) extension specifically designed to

integrate seamlessly with Ollama, a tool that allows users to run Large Language Models (LLMs) locally.<sup>46</sup> This extension brings Ollama's capabilities for local model execution directly into the VS Code environment, primarily focusing on providing efficient code autocompletion. A primary feature of Ollama Autocoder is its ability to execute LLMs locally based on the Ollama setup.<sup>46</sup> It supports switching between multiple models that are managed by Ollama, giving developers the flexibility to use different models for various coding tasks.<sup>46</sup> The extension is designed to provide low-latency responses, ensuring that code suggestions appear quickly as the developer types.<sup>46</sup> It also integrates web search functionality, allowing the local LLM to incorporate real-time information into its responses and suggestions, with intelligent synthesis of search results and accurate citation.<sup>46</sup> Ollama Autocoder includes an intelligent chat interface that provides streaming output of responses and even visualizes the model's thought process.<sup>46</sup> It also preserves the chat history for context in ongoing conversations.<sup>46</sup>

The extension offers flexible configuration options, allowing users to customize the server address for their Ollama instance, adjust performance modes, and configure various model parameters.<sup>46</sup> Using Ollama Autocoder is straightforward: in a text document, pressing the spacebar (or any character defined in the completion keys setting) will trigger the autocompletion feature, presenting an option to "Autocomplete with Ollama" or a preview of the first line of the suggestion.<sup>47</sup> Pressing Enter will then start the generation of the code. Alternatively, users can trigger the autocompletion via a command in the command palette or by setting a custom keybinding.<sup>47</sup> The extension streams the generated tokens directly to the cursor, and users can stop the generation early by pressing the "Cancel" button in the Ollama Autocoder notification or by typing something.<sup>47</sup> By default, Ollama Autocoder is configured to use the qwen2.5-coder:latest model, highlighting its focus on code-related tasks.<sup>47</sup> This extension provides a simple and effective way for VS Code users to leverage the power of locally run Ollama models for code autocompletion and chat assistance.

# Specialized Models for Code Generation and Understanding

While general-purpose Large Language Models (LLMs) have shown impressive capabilities in various natural language tasks, models specifically trained on code often exhibit superior performance when it comes to code generation, understanding, and related tasks. These specialized models are designed to grasp the nuances of programming languages, software development practices, and common coding patterns, making them invaluable tools for developers.

## Qwen2.5-Coder: Empowering Developers with State-of-the-Art Coding Abilities

Qwen2.5-Coder is a series of large language models developed by the Qwen team at Alibaba Cloud, specifically engineered for code-related tasks.<sup>48</sup> This series is presented as "Powerful,"

"Diverse," and "Practical," reflecting its aim to provide cutting-edge coding capabilities in a variety of contexts.<sup>48</sup>

The Qwen2.5-Coder series includes models of various sizes, ranging from 0.5 billion to 32 billion parameters, offering a spectrum of options to meet different computational needs.<sup>48</sup> Notably, the Qwen2.5-Coder-32B-Instruct model is claimed to be a state-of-the-art open-source code model, demonstrating coding abilities that are on par with those of GPT-40.<sup>48</sup> These models exhibit strong and comprehensive coding skills, along with solid general knowledge and mathematical reasoning abilities.<sup>48</sup> A significant feature of the Qwen2.5-Coder models is their support for a long context length of 128,000 tokens, enabling them to understand and generate longer sequences of code, which is crucial for handling complex and extensive codebases.<sup>48</sup> The series supports an impressive 92 coding languages, making it a versatile tool for developers working across different programming paradigms.<sup>48</sup> Furthermore, these code models retain the general and mathematical strengths of their base models, ensuring well-rounded performance.<sup>48</sup>

The Qwen2.5-Coder models are designed for practical applications in scenarios such as code assistants and Artifacts, indicating their potential for real-world use in development workflows.<sup>48</sup> Models with the "-Instruct" suffix are specifically designed for chatting and can be effectively used as code assistants, allowing developers to ask coding-related questions and receive helpful explanations and suggestions.<sup>48</sup> Base models, without the "-Instruct" suffix, are typically employed for code completion tasks, where the model continues code snippets based on the provided context.<sup>48</sup> For handling extremely long code inputs that might exceed the standard context length, Qwen2.5-Coder supports the YaRN technique.<sup>48</sup> It can also perform file-level code completion by filling in missing code segments within a given context, using special tokens to denote the prefix, suffix, and middle parts of the code.<sup>48</sup> Moreover, the model has the capability for repository-level code completion, where it can understand the relationships between different files within a project and complete code accordingly, using special tokens to indicate the repository structure and separate files.<sup>48</sup> With its claim of matching GPT-40's coding prowess and its versatility in handling various coding tasks, Qwen2.5-Coder stands out as a powerful open-source option for developers seeking advanced AI-powered code assistance.

## Honorable Mentions (Local Based)

In addition to the detailed reviews above, several other local-based AI tools warrant a brief mention for their contributions to the landscape:

- **Tabby** <sup>38</sup>: As discussed earlier, Tabby is an open-source, self-hosted AI coding assistant that provides an alternative to cloud-based solutions, emphasizing control and customization.
- **Privy** <sup>44</sup>: Also detailed previously, Privy is a VS Code extension that focuses on local AI autocomplete and chat using models like Ollama, prioritizing privacy.
- **Ollama Autocoder** <sup>46</sup>: Previously reviewed, this VS Code extension simplifies the integration of locally run Ollama models for code completion and chat.

- **Ollama** <sup>53</sup>: Ollama itself is a crucial tool that enables running various Large Language Models locally, providing a user-friendly interface and API for interacting with these models.<sup>55</sup> It supports a wide range of models, including Llama 2, Mistral, and others, and offers features like model management, customization, and hardware acceleration.<sup>55</sup>
- **Ilama cpp** <sup>59</sup>: This is a C++ inference library that aims to enable efficient LLM inference on a wide range of hardware, including CPUs, with minimal setup.<sup>61</sup> It supports various models and quantization techniques, making it possible to run powerful LLMs on consumer-grade hardware.<sup>59</sup>

These honorable mentions represent a selection of noteworthy local AI tools and frameworks that offer various functionalities for code assistance and LLM interaction, further highlighting the growing ecosystem of locally hosted AI solutions.

## Conclusion: Embracing the AI Tooling Era for Enhanced Productivity and Innovation

The exploration of AI tooling across search engines, local AI models, and IDE plugins reveals a dynamic and rapidly evolving landscape. Tools like Perplexity AI, Phind, and Gemini are redefining how we seek and process information, offering more intelligent and efficient alternatives to traditional search methods. The rise of local AI models such as Pieces, Goose, aichat, and Jan underscores a growing demand for privacy, control, and the ability to run powerful AI capabilities directly on personal hardware. Furthermore, the integration of AI into IDEs through plugins like ProxyAI, Tabby, Privy, and Ollama Autocoder is transforming the software development experience, providing developers with intelligent assistance that streamlines coding workflows and enhances productivity.

The initial positive sentiment towards AI tooling appears to be well-founded, given the significant advancements and the tangible benefits these tools offer across various domains. The trend towards AI-powered search is providing users with more direct, sourced, and contextually relevant answers, saving time and improving the quality of information retrieval. The development of local AI solutions is empowering users with greater privacy and control over their AI interactions, while also democratizing access to powerful language models. The deep integration of AI into IDEs is poised to revolutionize software development, making coding more efficient, less error-prone, and ultimately more productive.

Looking ahead, the evolution of AI tooling is expected to continue at a rapid pace, with new tools, features, and more powerful underlying models constantly emerging. Embracing these advancements will be crucial for individuals and organizations seeking to enhance their productivity, streamline their workflows, and foster innovation in an increasingly AI-driven world. The hype surrounding AI tooling appears justified by the current state of these technologies and their potential to significantly impact how we work and interact with information.

#### Works cited

- 1. What is Perplexity Al? How it Works, Key Features, Use Cases, & More | Guru, accessed May 4, 2025,
  - https://www.getguru.com/reference/what-is-perplexity-ai-and-how-to-use-it
- Perplexity Al Review: Features, Benefits, and Alternatives ClickUp, accessed May 4, 2025, <u>https://clickup.com/blog/perplexity-ai-review/</u>
- 3. How does Perplexity work? | Perplexity Help Center, accessed May 4, 2025, https://www.perplexity.ai/help-center/en/articles/10352895-how-does-perplexitywork
- 4. Introducing Perplexity Deep Research, accessed May 4, 2025, <u>https://www.perplexity.ai/hub/blog/introducing-perplexity-deep-research</u>
- 5. I replaced Google Search with Perplexity for a week: Here's why I liked and what I didn't, accessed May 4, 2025, https://www.androidpolice.com/replacied-google-search-wih-perplexity-for-a-w eek/
- 6. Perplexity, not Google, is now the best search engine The Register, accessed May 4, 2025,
  - https://www.theregister.com/2024/12/16/opinion\_column\_perplexity\_vs\_google/
- 7. The Latest on Perplexity Al vs. Google WebFX, accessed May 4, 2025, <u>https://www.webfx.com/blog/seo/perplexity-ai-vs-google/</u>
- 8. Perplexity vs. Google: Which Search Tool is Best in 2024? ClickUp, accessed May 4, 2025, <u>https://clickup.com/blog/perplexity-vs-google/</u>
- 9. How Perplexity could defeat Google (if it doesn't fumble) Command Al, accessed May 4, 2025, <u>https://www.command.ai/blog/perplexity-vs-google/</u>
- 10. I Tried Perplexity For a Week, And I Don't Think AI Search Engines Can Replace Google.. Yet | HackerNoon, accessed May 4, 2025, <u>https://hackernoon.com/i-tried-perplexity-for-a-week-and-i-dont-think-ai-searc</u> <u>h-engines-can-replace-google-yet</u>
- 11. What do you prefer: Perplexity or Google Search? : r/perplexity\_ai Reddit, accessed May 4, 2025, <u>https://www.reddit.com/r/perplexity\_ai/comments/1btcvxq/what\_do\_you\_prefer\_p</u> <u>erplexity\_or\_google\_search/</u>
- 12. Phind Review: Features, Pros, Cons, & Alternatives, accessed May 4, 2025, https://10web.io/ai-tools/phind/
- 13. Phind Features, Pricing, and Alternatives | Al Tools, accessed May 4, 2025, https://aitools.inc/tools/phind
- 14. Phind Review 2025 Features, Pricing & Deals ToolsForHumans.ai, accessed May 4, 2025, <u>https://www.toolsforhumans.ai/ai-tools/phind</u>
- 15. phind-70b model, accessed May 4, 2025, https://www.phind.com/search?cache=zl73ckwfaxfn7yl9uq35328w
- 16. Phind Al Answer Engine: Everything You Need To Know Fello Al, accessed May 4, 2025, <u>https://felloai.com/2024/12/phind-ai-answer-engine-everything-you-need-to-kno</u> w/
- 17. Phind 2: Al search with visual answers and multi-step reasoning | Hacker News, accessed May 4, 2025, <u>https://news.ycombinator.com/item?id=43039308</u>

- 18. Does any use phind and can give a comparison to perplexity? : r/perplexity\_ai -Reddit, accessed May 4, 2025, <u>https://www.reddit.com/r/perplexity\_ai/comments/1c9dpkk/does\_any\_use\_phind\_and\_can\_give\_a\_comparison\_to/</u>
- 19. The Gemini ecosystem Google Al, accessed May 4, 2025, https://ai.google/get-started/gemini-ecosystem/
- 20. How is Google Gemini Enhancing Search Capabilities? Thunder::Tech, accessed May 4, 2025, <u>https://www.thundertech.com/blog-news/how-is-google-gemini-enhancing-sear</u> <u>ch-capabilities</u>
- 21. Generative AI in Search: Let Google do the searching for you, accessed May 4, 2025,

https://blog.google/products/search/generative-ai-google-search-may-2024/

- 22. Al Tools for Business | Google Workspace, accessed May 4, 2025, https://workspace.google.com/solutions/ai/
- 23. Google One Al Premium Plan and Features, accessed May 4, 2025, https://one.google.com/about/ai-premium/
- 24. New Gemini app features, available to try at no cost, accessed May 4, 2025, <u>https://blog.google/products/gemini/new-gemini-app-features-march-2025/</u>
- 25. Gemini Advanced get access to Google's most capable AI models with Gemini 2.0, accessed May 4, 2025, <u>https://gemini.google/advanced/</u>
- 26. Introducing Gemini, your new personal Al assistant, accessed May 4, 2025, <u>https://gemini.google/assistant/</u>
- 27. Artificial Intelligence : Google Gemini Research Guides Syracuse University, accessed May 4, 2025,
  - https://researchguides.library.syr.edu/c.php?g=1341750&p=10240017 What you can do with your Gemini mobile app - Android - Google Help, acc
- 28. What you can do with your Gemini mobile app Android Google Help, accessed May 4, 2025, https://support.google.com/gemini/apswor/145796212bl-on&co-GENIE Platform%

https://support.google.com/gemini/answer/14579631?hl=en&co=GENIE.Platform% 3DAndroid

- 29. How is Google an amazing search engine but Gemini is so bad? Reddit, accessed May 4, 2025, <u>https://www.reddit.com/r/NoStupidQuestions/comments/1isbqh5/how\_is\_google\_</u> an amazing search engine but gemini/
- 30. I tried Google's new AI mode powered by Gemini, and it might be the end of Search as we know it | TechRadar, accessed May 4, 2025, <u>https://www.techradar.com/computing/artificial-intelligence/i-tried-googles-new-ai-mode-powered-by-gemini-and-it-might-be-the-end-of-search-as-we-know-it</u>
- 31. Do you use Gemini or Google search for finding things out? : r/Bard Reddit, accessed May 4, 2025, <u>https://www.reddit.com/r/Bard/comments/1amwm3a/do\_you\_use\_gemini\_or\_goo\_gle\_search\_for\_finding/</u>
- 32. Generative AI vs. Traditional Search: Technical Differences Matthew Edgar, accessed May 4, 2025,

https://www.matthewedgar.net/generative-ai-vs-traditional-search-technical-diff erences/

- 33. pieces-app/cli-agent: Pieces CLI for interacting with Pieces .. GitHub, accessed May 4, 2025, <u>https://github.com/pieces-app/cli-agent</u>
- 34. block/goose: an open source, extensible AI agent that goes ... GitHub, accessed May 4, 2025, <u>https://github.com/block/goose</u>
- 35. sigoden/aichat: All-in-one LLM CLI tool featuring Shell ... GitHub, accessed May 4, 2025, <u>https://github.com/sigoden/aichat</u>
- 36. menloresearch/jan: Jan is an open source alternative to ... GitHub, accessed May 4, 2025, <u>https://github.com/menloresearch/jan</u>
- 37. carlrobertoh/ProxyAl: The leading open-source Al copilot ... GitHub, accessed May 4, 2025, <u>https://github.com/carlrobertoh/ProxyAl</u>
- 38. What is Tabby? Features & Getting Started Deepchecks, accessed May 4, 2025, <u>https://www.deepchecks.com/llm-tools/tabby/</u>
- 39. Tabby Free And Opensource Copilot Alternative With Code Example, accessed May 4, 2025, https://withcodeexample.com/tabby-free-and-opensource-copilot-alternative/
- 40. About the Docs Tabby, accessed May 4, 2025, https://tabby.tabbyml.com/docs/welcome/
- 41. Tabby Opensource, self-hosted AI coding assistant, accessed May 4, 2025, https://www.tabbyml.com/
- 42. What is tabby? Features & getting started for developers and SMBs BytePlus, accessed May 4, 2025, <u>https://www.byteplus.com/en/topic/555358</u>
- 43. What is Tabby AI? The ultimate guide for developers and SMBs BytePlus, accessed May 4, 2025, <u>https://www.byteplus.com/en/topic/555348</u>
- 44. Privy AI coding Autocomplete and chat that runs locally. Visual Studio Marketplace, accessed May 4, 2025, https://marketplace.visualstudio.com/items?itemName=privy.privy-vscode
- 45. Privy AI coding Autocomplete and chat that runs locally. Open VSX Registry, accessed May 4, 2025, <u>https://open-vsx.org/extension/Privy/privy-vscode/reviews</u>
- 46. VSCode Ollama Visual Studio Marketplace, accessed May 4, 2025, https://marketplace.visualstudio.com/items?itemName=warm3snow.vscode-olla ma
- 47. Ollama Autocoder Visual Studio Marketplace, accessed May 4, 2025, <u>https://marketplace.visualstudio.com/items?itemName=10nates.ollama-autocode</u> <u>r</u>
- 48. GitHub QwenLM/Qwen2.5-Coder, accessed May 4, 2025, https://github.com/QwenLM/Qwen2.5-Coder
- 49. Buy now, pay later How it works Tabby, accessed May 4, 2025, <u>https://tabby.ai/pay-later</u>
- 50. A Guide to Self-Hosted LLM Coding Assistants Semaphore, accessed May 4, 2025, <u>https://semaphoreci.com/blog/selfhosted-llm-coding-assistants</u>
- 51. autocompletion does not work in VSCode · Issue #37 · srikanth235/privy GitHub, accessed May 4, 2025, <u>https://github.com/srikanth235/privy/issues/37</u>
- 52. What VS Code Extension works best with Ollama models? Reddit, accessed May

4, 2025,

https://www.reddit.com/r/ollama/comments/1dvb2e1/what\_vs\_code\_extension\_w orks\_best\_with\_ollama/

- 53. Ollama Features | Elest.io, accessed May 4, 2025, https://elest.io/open-source/ollama/resources/software-features
- 54. PixelLlama 0.95b New features and improves : r/ollama Reddit, accessed May 4, 2025,

https://www.reddit.com/r/ollama/comments/1gvr55k/pixelllama\_095b\_new\_featur es\_and\_improves/

- 55. What is Ollama and how to use it: a quick guide [part 1] Geshan Manandhar, accessed May 4, 2025, <u>https://geshan.com.np/blog/2025/02/what-is-ollama/</u>
- 56. What is Ollama? Understanding how it works, main features and models -Hostinger, accessed May 4, 2025, https://www.hostinger.com/tutorials/what-is-ollama
- 57. ollama/ollama: Get up and running with Llama 3.3, DeepSeek-R1, Phi-4, Gemma 3, Mistral Small 3.1 and other large language models. - GitHub, accessed May 4, 2025, <u>https://github.com/ollama/ollama</u>
- 58. Blog · Ollama, accessed May 4, 2025, https://ollama.com/blog
- 59. Llama CPP Tutorial: A Basic Guide And Program For Efficient LLM Inference And Models, accessed May 4, 2025, <u>https://pwskills.com/blog/llama-cpp/</u>
- 60. llama.cpp guide Running LLMs locally, on any hardware, from scratch ::, accessed May 4, 2025, <u>https://steelph0enix.github.io/posts/llama-cpp-guide/</u>
- 61. ggml-org/llama.cpp: LLM inference in C/C++ GitHub, accessed May 4, 2025, https://github.com/ggml-org/llama.cpp
- 62. Llama.cpp Tutorial: A Complete Guide to Efficient LLM Inference and Implementation, accessed May 4, 2025, https://www.datacamp.com/tutorial/llama-cpp-tutorial
- 63. Ilama.cpp: The Ultimate Guide to Efficient LLM Inference and Applications -PyImageSearch, accessed May 4, 2025, <u>https://pyimagesearch.com/2024/08/26/Ilama-cpp-the-ultimate-guide-to-efficien</u> <u>t-Ilm-inference-and-applications/</u>
- 64. What is Llama.cpp? Key Features & Getting Started Deepchecks, accessed May 4, 2025, <u>https://www.deepchecks.com/llm-tools/llama-cpp/</u>
- 65. Feature matrix · ggml-org/llama.cpp Wiki GitHub, accessed May 4, 2025, <u>https://github.com/ggml-org/llama.cpp/wiki/Feature-matrix</u>