# The OpenHands Software Agent SDK: A Composable and Extensible Foundation for Production Agents

**Xingyao Wang**, **Simon Rosenberg**, **Juan Michelini**, **Calvin Smith**, **Hoang Tran**, **Engel Nyst**, **Rohit Malhotra**, **Xuhui Zhou**, **Valerie Chen**, **Robert Brennan**, **Graham Neubig**

{xingyao, graham}@openhands.dev

Agents are now used widely in the process of software development, but building production-ready software engineering agents is a complex task. Deploying software agents effectively requires flexibility in implementation and experimentation, reliable and secure execution, and interfaces for users to interact with agents. In this paper, we present the **OpenHands Software Agent SDK**, a toolkit for implementing software development agents that satisfy these desiderata. This toolkit is a complete architectural redesign of the agent components of the popular **OpenHands** framework for software development agents, which has 64k+ GitHub stars.

To achieve flexibility, we design a *simple interface for implementing agents* that requires only a few lines of code in the default case, but is easily extensible to more complex full-featured agents with features such as custom tools, memory management, and more. For security and reliability, it delivers *seamless local-to-remote execution portability*, integrated REST/WebSocket services. For interaction with human users, it can *connect directly to a variety of interfaces*, such as visual workspaces (VS Code, VNC, browser), command-line interfaces, and APIs. Compared with existing SDKs from OpenAI, Claude and Google, OpenHands uniquely integrates native sandboxed execution, lifecycle control, model-agnostic multi-LLM routing, and built-in security analysis. Empirical results on SWE-Bench Verified and GAIA benchmarks demonstrate strong performance. Put together, these elements allow the OpenHands Software Agent SDK to provide a practical foundation for prototyping, unlocking new classes of custom applications, *and* reliably deploying agents at scale.
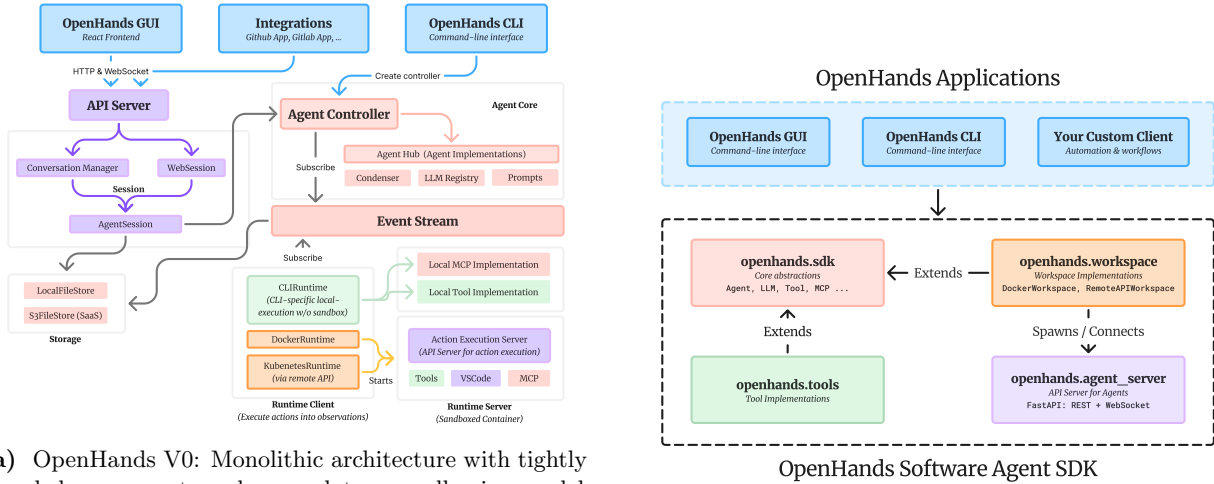
**Software Agent SDK:** https://github.com/OpenHands/software-agent-sdk

**Benchmarks:** https://github.com/OpenHands/benchmarks

## 1 Introduction

In software engineering, AI agents have evolved from assistive tools (GitHub, 2021; Cursor Team, 2024) to autonomous systems capable of hours-long async execution on complex tasks (e.g., Devin (Cognition AI, 2024), Claude Code (Anthropic, 2025b), OpenHands (Wang et al., 2025)). Reliably deploying autonomous agents in production requires *system foundations*—including durable state management, safe execution in a sandbox, and consistent behavior across environments ranging from local to containerized cloud deployments—that were unnecessary for earlier assistive tools that rely on largely local user-driven workflows.

OpenHands (Wang et al., 2025) demonstrated that open-source software agents can achieve broad adoption and contributions, reaching over 64k GitHub stars and hundreds of contributors in just 18 months (Neubig, 2025). As the project scaled, the original monolithic architecture – what we refer to as OpenHands V0 – exposed growing architectural tensions. The early design was driven by the need for fast prototyping and iteration, combining agent logic, evaluation, and applications in a single codebase. Over time, this approach led to rigid sandboxing assumptions, sprawling mutable configuration, and tight coupling between research

**(a)** OpenHands V0: Monolithic architecture with tightly coupled components and a mandatory sandboxing model. The design assumed all executions occur in a sandbox, making later support for local execution workflows (CLI) cumbersome — requiring special-case handling in the CLI runtime and duplicated local implementations of MCP and tools.



**(b)** OpenHands V1: Modular SDK architecture with four decoupled packages. Applications consume software agent SDK; `tool` and `workspace` packages extend `sdk`'s abstraction with actual implementations; the `workspace` package spawns and connects to `agent server`.

**Figure 1** Architecture Evolution from OpenHands V0 to V1. V0's monolithic, sandbox-centric design tightly coupled components and required duplicated local implementations. V1 refactors this into a modular SDK with clear boundaries, opt-in sandboxing, and reusable agent, tool, and workspace packages. The color of each component in V0 roughly corresponds to its modular counterpart in V1.

and production, eventually making a full architectural redesign unavoidable.

Guided by these lessons, **OpenHands V1** introduces a new architecture grounded in four design principles that directly address the limitations of V0:

- **Optional isolation.** The agent runs locally by default but can switch to a sandboxed environment when additional safety or resource control is required.
- **Stateless by default, one source of truth for state.** All components—agents, tools, LLMs, etc—are immutable and validated at construction, while a single conversation state object records all mutable context. This design ensures reliable recovery of agent sessions.
- **Strict separation of concerns.** The agent core is decoupled from applications so that downstream systems (CLI, Web UI, GitHub App) use it as a shared library rather than duplicating logic.
- **Two-layer composability.** Developers can compose independent deployment packages (SDK, Tools, Workspace, Server) and extend the SDK safely by adding or replacing typed components such as tools or agents.

Building on these principles, we introduce **OpenHands V1**, a complete software-agent ecosystem—including CLI and GUI applications—built on a shared foundation: the **OpenHands Software Agent SDK** (Fig. 1b). The SDK defines an *event-sourced state model* with deterministic replay, an *immutable configuration* for agents, and a *typed tool system* with MCP integration. Its *workspace abstraction* enables the same agent to run locally for prototyping or remotely in secure, containerized environments with minimal code changes. Unlike prior library-only SDKs (Anthropic, 2025a; OpenAI, 2024), OpenHands includes a *built-in REST/WebSocket server* for remote execution and a *suite of interactive workspace interfaces*—a browser-based VSCode IDE, VNC desktop, and persistent Chromium browser—for human inspection and control.

We systematically compare 31 features of our SDK with those of the OpenAI Agents SDK, Claude Agent SDK,

```python
from openhands.sdk import LLM, Conversation
from openhands.tools.preset.default import get_default_agent
llm = LLM(model="openhands/claude-sonnet-4-5-20250929", api_key="...")
agent = get_default_agent(llm=llm)
conversation = Conversation(agent=agent, workspace="/path/to/project")
conversation.send_message("Write 3 facts about this project into FACTS.txt.")
conversation.run()
```

**Figure 2** Minimal example of OpenHands Software Agent SDK.

and Google ADK, and find that, although 15 features are shared with at least one of them, our SDK uniquely combines 16 additional features, including native remote execution, a production server with sandboxing, and model-agnostic multi-LLM routing across 100+ providers. Our SDK further adds a security analyzer for agent actions, flexible lifecycle control (pause/resume, sub-agent delegation, history restore, etc), and built-in QA instrumentation (unit tests, LLM-based integration tests, and evaluation benchmarks) for production reliability.

Across multiple LLM backends, our SDK achieves strong results on SWE-Bench Verified and GAIA benchmarks, demonstrating both state-of-the-art performance and the generality of the architecture to operate consistently across diverse model providers.

The OpenHands Software Agent SDK is fully open-sourced under the MIT License at https://github.com/OpenHands/software-agent-sdk.

# 2 Preliminaries

## 2.1 Agent Design

Modern AI agents can be viewed as systems that perceive, reason, and act within an environment to accomplish goals over time (Russell and Norvig, 2009). In recent work (Wang et al., 2023), agent design typically centers on these interconnected components:

- **Environment and Observations.** The environment defines the external context in which the agent operates, providing structured interfaces (e.g., APIs, files, user interfaces) from which it receives observations. These percepts—such as text, logs, or visual data—inform the agent's situational awareness and ground its decisions in the evolving task state.
- **State and Memory.** An agent's state encodes what it currently knows about the world and itself, often represented as a history of past interactions, retrieved information, and internal reasoning.
- **Actions and Tools.** Actions represent how the agent influences its environment—issuing commands, executing code, or invoking APIs—each realized through the use of tools that implement these capabilities.

## 2.2 Software Agents

In software engineering contexts, these agentic components manifest concretely in systems like Open-Hands (Wang et al., 2025), which emulate a developer's end-to-end workflow. The environment is a sandboxed workspace exposing a filesystem, terminal, and web interface, allowing the agent to observe program output and test feedback. Its state and memory are represented through an event log that records commands, edits, and results, providing persistent context across actions. The agent acts by editing files, running tests, or invoking structured tools like web browser to interact with its environment.

# 3 Challenges and Design Principles

As OpenHands evolved, it has accumulated significant architectural complexity: core agent logic, multiple CLIs, a web server for the GUI, frontend code, and diverse runtime providers (e.g., local Docker, Kubernetes in production) all coexist in a single codebase. Much of this architecture dated back to the project's first weeks, before the introduction of LLM structured tool use or the Model Context Protocol (MCP; MCP Team 2025). This section revisits the key architectural tensions observed in OpenHands V0 and distills them into four design principles guiding V1.

## 3.1 Universal Sandboxing vs. Local Flexibility

**V0 Challenges.** V0 was built on the assumption that all tool calls should run inside sandboxed Docker containers for safety and reproducibility.

However, doing so introduced multiple layers of friction as each conversation spanned two independent processes (agent and sandbox) with potentially divergent states. Since the sandbox for tools might crash while the agent continued (or vice versa), leading to corrupted sessions. In multi-tenant deployments, resource-heavy actions from one user, such as mass webpage visits generating oversized screenshots, could exhaust container resources and crash other agents sharing the same application.

When the need arose to support local execution workflows—for example, running tools or MCP clients directly on the host machine—we had to add exceptions and bypass layers that duplicated existing logic. As shown in Fig. 1a, this led to redundant local versions of MCP and tool implementations, diverging from the original sandbox-based code path. These ad hoc extensions made the architecture increasingly brittle and poorly aligned with MCP, which assumes agents can execute actions locally with direct access to credentials, files, and IDEs.

> **V1 Design Principle**
>
> **Sandboxing should be opt-in, not universal.** V1 unifies agent and tool execution in a single process by default, aligning with MCP's assumptions. When isolation is needed, the same stack can be containerized transparently. Sandboxing becomes opt-in—preserving flexibility without sacrificing safety.

## 3.2 Mutable Configuration vs. Deterministic State

**V0 Challenges.** V0's configuration system mixed different domains — deployment, agent behavior, LLM routing, sandbox settings, and UI options — within overlapping layers, creating hidden dependencies and complex override logic. The deeper issue was structural: configuration was split across multiple parallel hierarchies (CLI/headless, Web UI, GitHub App, and SaaS), each evolving its own precedence rules and assumptions. Different entry points were added incrementally, each patching configuration values in place, so two runs with identical parameters could still diverge subtly. Later attempts to unify these systems only added more override layers and inconsistencies, while the web configuration remained rigid due to its coupling with ORM models and database schemas. These issues led to severe sprawl — 140+ fields, 15 classes, and 2.8K lines of configuration code — a brittle system where small changes often cascaded into unrelated failures.

> **V1 Design Principle**
>
> **Stateless by Default, One Source of Truth for State.** V1 treats all agents and their components—tools, LLMs, etc—as immutable and serializable Pydantic models validated at construction. The only mutable entity is the *conversation state*, which is a single, well-defined source of truth that tracks ongoing execution. This design isolates change to one place, enabling deterministic replay, strong consistency, and stable long-term recovery.

## 3.3 Monorepo vs. Modular SDK

**V0 Challenges.** OpenHands V0 was implemented as a single monolithic repository that combined the agent core, evaluation suite, and applications (frontend, backend, CLI) into one codebase. This monorepo simplified early development but gradually blurred boundaries between the agent and its downstream applications. The same repository powered the CLI, web interface, and GitHub integrations, each introducing divergent logic and configuration overrides. Over time, the agent core absorbed application-specific branches, benchmark dependencies, and environment-specific hacks, making the system brittle and slow to evolve.

As OpenHands gained popularity in the academic community, many benchmarks were contributed, making the agent more general but also introducing heavy dependencies and frequent version conflicts. These leaked into the main application due to mono-repo design, making deployments heavyweight and fragile.

> **V1 Design Principle**
>
> **Maintain strict separation of concerns.** V1 isolates the agent core into software engineering SDK as described in this paper. Applications integrate via SDK APIs, allowing research to evolve independently from applications.

## 3.4 Monolith Logic vs. Extensible Architecture

**V0 Challenges.** As V0 lacked clear boundaries between the agent core and its applications, adding new behaviors in V0 often required editing the core logic or branching for specific entry points, limiting experimentation and maintainability. The system lacked a structured notion of composability and extensibility, forcing ad hoc hacks for even small changes.

> **V1 Design Principle**
>
> **Everything should be composable and safe to extend.** V1 makes composability a first-class design goal at two levels. At the *deployment level*, its four modular packages—SDK, Tools, Workspace, and Agent Server—combine flexibly to support local, hosted, or containerized execution. At the *capability level*, the SDK exposes a typed component model—tools, LLMs, contexts, etc—so developers can extend or reconfigure agents declaratively without touching the core.

# 4 Architecture

OpenHands V1's architecture emerged from both the operational challenges we encountered in V0 analyzed in §3. OpenHands V1 is a *broader ecosystem* with applications like CLI and GUI. Its foundation is the **OpenHands Software Agent SDK**—a standalone developer framework comprising nine interlocking components: Event-Sourced State Management (§4.2), LLM (§4.3), Tool System (§4.4), Agent (§4.5), Context Window Management (§4.6), Local Conversation (§4.7), Secret Registry (§4.8), Security and Confirmation (§4.9), and deployment architecture (§4.10) with local and remote workspace support (§4.10). This section focuses on describing the *SDK architecture* rather than the full application and explains how these components collectively support both local and production deployments.

## 4.1 Modular Four-Package Design

The OpenHands Software Agent SDK is organized into four Python packages that compose together based on deployment needs. Fig. 1b shows how these packages interact:

- `openhands.sdk`: Core abstractions (Agent, Conversation, LLM, Tool, MCP, etc) and the event system.
- `openhands.tools`: Concrete tool implementations based on abstractions defined in `openhands.sdk`.

- openhands.workspace: Execution environments (e.g., Docker, hosted API) that extend SDK base classes.
- openhands.agent_server: A web server exposing REST/WebSocket APIs for remote execution.

The separation addresses key production concerns that slowed down development, QA, and release in OpenHands V0: (1) sdk stays lightweight for diverse integration scenarios; (2) tools isolates slow-running tool tests from core SDK changes, speeding up development; (3) workspace provides optional sandboxing implementations without bloating the core; and (4) agent_server offers a generic API server usable with or without containers. This modularity enables independent testing, selective dependency management, and incremental release cycles—critical for production deployments where monolithic repositories create QA bottlenecks.

## 4.2   Event-Sourced State Management

At V1's core lies an event-sourcing pattern treating all interactions as immutable events appended to a log.

**Event Hierarchy.** The event system uses a multi-level hierarchy shown in Table 1. At the base, Event provides immutable structure (ID, timestamp, source) with type-safe serialization via discriminated unions (Pydantic Team, 2025). LLMConvertibleEvent adds to_llm_message() for converting events into LLM format. The action-observation loop uses ActionEvent for tool calls and ObservationBaseEvent subclasses for results – all these are subclasses of LLMConvertibleEvent. Internal events (condensation, state updates, etc) inherit directly from Event for bookkeeping without LLM exposure.

**Table 1**   Event hierarchy organized by parent-child relationships. LLM-convertible events can be sent to the LLM, while internal events handle state management and control flow.

| Parent Class | Child Classes & Purpose |
|---|---|
| Event | *Base:* Immutable structure, type-safe serialization |
| LLMConvertibleEvent<br>↪ Event | *LLM-Convertible Events (visible to LLM)*<br>MessageEvent: User/assistant text messages<br>ActionEvent: Agent tool calls with thought & reasoning<br>SystemPromptEvent: System prompt with tool schemas<br>CondensationSummaryEvent: Summary of forgotten events<br>ObservationBaseEvent: Base for tool responses |
| ObservationBaseEvent<br>↪ LLMConvertibleEvent | ObservationEvent: Successful tool execution results<br>UserRejectObservation: User rejected action<br>AgentErrorEvent: Agent/scaffold errors |
| Event | *Internal Events (not visible to LLM)*<br>ConversationStateUpdateEvent: State field changes<br>CondensationRequest: Trigger history compression<br>Condensation: Compression result with forgotten events<br>PauseEvent: User-requested pause |

**ConversationState: Single Source of Truth.** By design, components like Agent, Tool, and LLM are immutable and serializable—all changing variables live in ConversationState, making it the *only* stateful component. This class maintains two types of state: (1) mutable metadata fields (e.g., agent_status, stats, confirmation_policy) stored directly in the Pydantic model, and (2) an append-only EventLog recording all agent interactions. A FIFO lock ensures thread-safe updates through a two-path pattern: state-only updates for metadata changes, and event-based updates that append to the log.

When persistence is configured, ConversationState selectively writes changes to disk. Metadata fields serialize to a single base_state.json file on each modification, while EventStore persists events as individual JSON files to the corresponding directory. This dual-path design enables efficient incremental persistence—only new events write to disk, avoiding rewrites of large histories. Conversations resume by loading base_state.json

```python
class RouterLLM(LLM):
    llms_for_routing: dict[str, LLM]  # Available models

    @abstractmethod
    def select_llm(self, messages: list[Message]) -> str:
        """Return key of LLM to use from llms_for_routing."""

    def completion(
        self, messages: list[Message], ... **kwargs,
    ) -> LLMResponse:
        # Select appropriate LLM
        selected_model = self.select_llm(messages)
        self.active_llm = self.llms_for_routing[selected_model]
        return self.active_llm.completion(...)

class MultimodalRouter(RouterLLM):
    def select_llm(self, messages: list[Message]) -> str:
        has_images = any(m.contains_image for m in messages)
        return "primary" if has_images else "secondary"

# Usage: route text to cheaper model, images to multimodal model
router = MultimodalRouter(llms_for_routing={
    "primary": LLM(model="claude-sonnet-4-5"),
    "secondary": LLM(model="devstral-small")})
agent = Agent(llm=router, tools=tools)
```

**Figure 3** Multi-LLM routing example. Routers inherit from LLM to maintain a unified interface while delegating to selected models.

and replaying events from the directory, with agents automatically detecting incomplete conversations and continuing from the last processed event.

## 4.3 LLM Abstraction Layer

The LLM class provides a unified interface to language models. Through LiteLLM, it supports 100+ providers with two APIs: the standard Chat Completions API for broad compatibility and the newer OpenAI Responses API for latest reasoning models.

**Native Support for Reasoning / Extended Thinking.** The SDK captures and processes advanced native reasoning fields from frontier models, such as ThinkingBlock for Anthropic's extended thinking, and ReasoningItemModel for OpenAI's reasoning. The SDK supports the OpenAI Responses API transparently for the agent, enabling client developers to use the agent with advanced reasoning models like GPT-5-Codex that are only available on the recently released Responses API.

**Built-in Support for Non-function-calling Models.** For models without native function calling support, the SDK implements a NonNativeToolCallingMixin, which converts tool schemas to text-based prompt instructions and parses tool calls from model outputs using structured prompts and regex-based extraction. This enables models that do not support function calling to be used for agentic tasks, dramatically expanding the set of usable models.

**Multi-LLM Routing Support.** SDK features RouterLLM, a subclass of LLM that enables the agent to use different models for different LLM requests. Custom implementations can extend RouterLLM and implement select_llm() to choose a different model based on different LLM inputs. Please refer to Fig. 3 for a pseudo-code example.

## 4.4 Tool System

The V1 tool system provides a type-safe and extensible framework grounded in an **Action–Execution–Observation** pattern. As illustrated in Fig. 4, tool usage follows a strict contract: the LLM proposes JSON tool calls, which are validated and parsed into Action; these are executed and the results are returned as Observation.

This abstraction unifies how the SDK supports both custom tools and MCP tools, providing a single standard interface for defining, invoking, and managing tools.

**Action–Execution–Observation Pattern.** Each tool is defined by three well-separated components:

- **Action** — Specifies the input schema for a tool call. LLM-generated arguments are validated against a Pydantic model before execution, ensuring type safety and preventing malformed requests.

- **Execution** — Implements the tool's actual logic via the `ToolExecutor`, which receives a validated `Action` and performs the underlying execution.

- **Observation** — Captures the output of the execution, defining a structured return schema, and converting results (or errors) into a LLM-compatible format.

```python
class Schema(BaseModel):
    @classmethod
    def to_mcp_schema(cls) -> dict[str, Any]
    @classmethod
    def from_mcp_schema(cls: type[S], model_name: str, schema: dict[str, Any]) -> type["S"]

class Action(Schema):
    """Tool input with schema validation."""
    def visualize(self) -> str  # For UI display

class Observation(Schema):
    """Tool output with LLM conversion."""
    def to_llm_content(self) -> str  # For LLM context

class ToolDefinition[ActionT, ObservationT](BaseModel):
    action_type: type[Action]
    observation_type: type[Observation]
    executor: ToolExecutor
    def to_mcp_tool(self, ...) -> dict
    def to_openai_tool(self) -> dict  # ChatCompletions API format
    def to_responses_tool(self) -> dict  # Responses API format

class ToolExecutor[ActionT, ObservationT](ABC):
    """Executor function type for a Tool."""
    def __call__(self, action: ActionT) -> ObservationT
```

**Figure 4** Tool system structure. Actions validate inputs, executors run logic, and observations format outputs for LLMs.

**MCP Integration.** The same abstraction additionally enables seamless support for the MCP. MCP tools are treated as first-class SDK tools: their JSON Schemas are automatically translated into `Action` models, and their results are surfaced as structured `Observation`. `MCPToolDefinition` extends the standard `ToolDefinition` interface, while `MCPToolExecutor` delegates execution to FastMCP's `MCPClient`, which manages server communication and transport details. As a result, external MCP tools behave identically to native tools—validated on input, type-safe at runtime, and serialized for LLM consumption—highlighting the flexibility and extensibility of the core tool abstraction.

**Tool Registry and Distributed Execution.** SDK supports distributed agent architectures by decoupling *tool specifications* from *implementations* using a registry-based resolution mechanism. Because Python executors are non-serializable, tools are represented as lightweight `Tool` specification objects containing only a registered name and JSON-serializable parameters. Through `register_tool(name, ToolDefinition)`, each identifier is bound to a resolver that reconstructs the full definition—including its executor—at runtime based on conversation context. This allows tool specs to *cross process or network boundaries* as pure JSON, enables lazy instantiation with environment-specific state (e.g., workspace paths), and ensures a uniform interface for local and remote execution.

## 4.5   Agent: Stateless Event Processor

The agent abstraction separates *configuration* from *execution state*. Agents are defined as stateless, immutable specifications, including LLM settings, tool specifications, security policies, and agent content, that can be serialized and transmitted across process boundaries.

**Event-Driven Execution.** Agents execute through an event-driven loop that processes conversation state step-by-step. Rather than directly returning results, agents emit structured events (e.g., messages, actions, observations) through callbacks, i.e. `on_event(event: Event) -> None`, separating event generation from execution control. This design enables: (1) **security interleaving**—actions can be reviewed or blocked before execution based on risk analysis (§4.9); (2) **incremental execution**—the agent advances one step at a time, supporting pause/resume, recovery from context overflows, and condensation for long conversations; and (3) **event streaming**—intermediate results (e.g., observations and reasoning traces) are emitted in real time for UI updates and monitoring.

**Customizing Agent Context with Skills and Prompts.** `AgentContext` centralizes all inputs that shape LLM behavior, including prefixes/suffixes for system/user messages and user-defined `Skill` objects. Skills can be defined programmatically or loaded from markdown files (e.g., `.openhands/skills/`, or compatible formats like `.cursorrules`, `agents.md`). Each skill may always be active (`trigger=None`) to persistently augment the system prompt, or conditionally activated via keyword matching based on user input; skills may also include MCP tools. This design enables rich contextual and behavioral customization without modifying agent logic.

**Sub-Agent Delegation.** The SDK supports hierarchical agent coordination through a delegation tool that demonstrates the extensibility of the tool abstraction. Sub-agents operate as independent conversations that inherit the parent's model configuration and workspace context, enabling structured parallelism and isolation without any changes to the core SDK. The current implementation provides blocking parallel execution, implemented as a standard tool in the `openhands.tools` package, where the parent agent spawns and monitors sub-agents until all tasks complete. This pattern exemplifies how complex coordination behaviors—such as asynchronous delegation, dynamic scheduling, or fault-tolerant recovery—can be implemented entirely as user-defined tools, reinforcing the SDK's design principle for extensibility that advanced agent orchestration requires no modification to the core framework.

## 4.6   Context Window Management

To ensure the ever-growing history fits inside the LLM's context, the `Condenser` system drops events and replaces them with summaries whenever the history grows too large.

The results of any given condensation are stored in the event log as a `CondensationEvent`. Before sending the event history to the LLM, the agent *applies* these condensation events by removing forgotten events and inserting summaries. This strategy lets the SDK preserve the entirety of the event log, regardless of condensation, while also keeping the condenser implementations stateless.

Condensers enable long-running conversations and, as a mechanism for constraining the tokens consumed per step, can reduce the overall cost of a conversation. `LLMSummarizingCondenser` (the default condenser) has been shown to reduce API costs by up to 2× with no degradation in agent performance (Smith, 2025).

## 4.7   Local Conversation

`LocalConversation` provides the simplest and most direct execution mode of the SDK, designed for rapid iteration and debuggability. It runs the full agent loop—LLM calls, tool invocation, event callbacks, and state updates—entirely in-process without network or container overhead. The core API offers intuitive control: developers can initialize a conversation with `Conversation(agent, workspace)`, send inputs using `send_message()`, start execution via `run()`, pause and resume with `pause()` and `run()` respectively, inspect

results through event callbacks or direct access of `.state` attribute. Pausing automatically persists state and emits a `PauseEvent`, allowing agents to resume from the same point later.

## 4.8   Secret Registry

`SecretRegistry` provides secure, late-bound, and remotely manageable credentials for tool execution. Each conversation maintains its own instance, ensuring strict per-session isolation. Tools access secrets only at execution time, and all secret values appearing in outputs are masked to prevent leakage. For example, the Bash Tool scans commands for secret keys, exports the referenced ones as environment variables, and replaces their occurrences in results with a constant mask (`<secret-hidden>`). Secrets may be static values or callables (e.g., token refreshers) and can originate from local stores or HTTP-based sources with secret-aware header handling. New custom sources can be easily added to integrate with external secret stores. All secrets are redacted during serialization and can be encrypted with a configurable cipher. They can also be updated mid-conversation—locally or through the agent server API—supporting live rotation without restarting agents.

## 4.9   Security and Confirmation

AI agents can perform risky actions, making security especially important when they run on a user's machine. To address this, SDK treats security as a first-class concern within the agent's control loop. Two abstractions form the core of this design: the `SecurityAnalyzer`, which rates each tool call as low, medium, high, or unknown risk, and the `ConfirmationPolicy`, which determines whether user approval is required before execution based on the action's details and assessed risk.

When approval is required, the agent pauses in a special `WAITING_FOR_CONFIRMATION` state until the user explicitly approves or rejects the action. Upon rejection, it may retry with safer alternatives. The confirmation policy can be updated dynamically during a session, enabling adaptive trust—such as relaxing restrictions for safe, read-only operations like grep.

This architecture separates risk assessment from enforcement, allowing developers define custom `SecurityAnalyzer` and `ConfirmationPolicy` without touching tool executors or core logic. The SDK includes a built-in pair: `LLMSecurityAnalyzer`, which appends a `security_risk` field to tool calls, and `ConfirmRisky` policy, which blocks actions exceeding a configurable risk threshold (default: high).

## 4.10   Deployment Architecture: Local to Remote

OpenHands Agent SDK's key innovation is seamless transition from local prototyping to remote production with minimal code changes. This is enabled by the agent API server and the abstraction of `Conversation` and `Workspace`.

**Conversation Factory: Local & Remote Conversation.** The `Conversation` class serves as a factory entry point that abstracts over local and remote execution. When instantiated with a string path or `LocalWorkspace`, it returns a `LocalConversation` that executes the full agent loop in-process, directly invoking tools and updating state synchronously (§4.7). When provided a `RemoteWorkspace`, the same call transparently constructs a `RemoteConversation`, which serializes the agent configuration and delegates execution to an agent server over HTTP and WebSocket. Both implementations share an identical API allowing seamless migration from local prototyping to containerized multi-user deployments without code changes. This factory pattern encapsulates all environment-specific logic behind a unified interface, achieving the SDK's goal of "local-first, deploy-anywhere" development: the same agent code that runs interactively in a notebook can scale to distributed production simply by swapping the workspace type.

**Agent Server.** The `agent_server` module implements an API server for remote agent execution (Fig. 6). It exposes REST endpoints for conversation control (e.g., `POST /conversations`, `GET /conversations/id`) and

```
1   --- local.py
2   +++ remote.py
3   @@
4    from openhands.sdk import LLM, Conversation
5    from openhands.sdk.preset.default import get_default_agent
6    llm = LLM(model="anthropic/claude-sonnet-4.1", ...)
7    agent = get_default_agent(llm=llm)
8   -conversation = Conversation(agent=agent)
9   -conversation.send_message("Create hello.py")
10  -conversation.run()
11  +from openhands.workspace import DockerWorkspace
12  +with DockerWorkspace(...) as workspace:
13  +  conversation = Conversation(agent=agent, workspace=workspace)
14  +  conversation.send_message("Create hello.py")
15  +  conversation.run()
```

**Figure 5** Local-to-remote transition requires only importing and instantiating `DockerWorkspace`. All other code (agent configuration, LLM setup, message handling) remains unchanged.

WebSocket for event streaming. When a `RemoteConversation` starts, it serializes agent configuration—including LLM settings, tools, and context—into JSON and submits it to `/conversations`. The server reconstructs the agent, launches a local execution loop, and streams structured events back in real time, enabling responsive UIs without polling overhead.
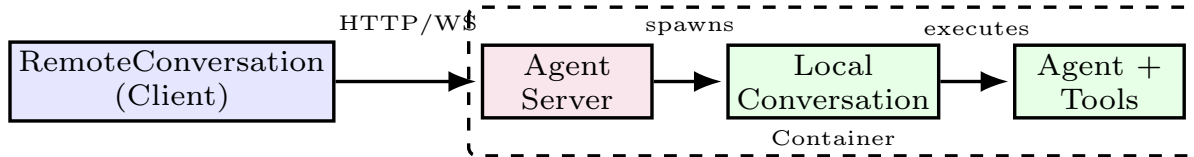


**Figure 6** Agent server architecture. Client serializes agent configuration via HTTP; the server executes using SDK components inside the container and streams events via WebSocket.

To support scalable and isolated execution, we provide official Docker images that bundle the full agent-server stack—including the API server, VSCode Web, VNC desktop, and Chromium browser. Each agent instance runs in an independent container with a dedicated file system, environment, and resource. This containerized design simplifies deployment and enables SaaS-style multi-tenancy while preserving workspace isolation.

**Workspace.** The `BaseWorkspace` abstract class enables sandboxes for agents. **Local Workspace** executes in-process against the host filesystem and shell; it is effectively a thin, no-op wrapper that forwards file/command/git operations directly, enabling fast prototyping without network hops. **Remote Workspace** preserves the same interface but delegates all operations over HTTP to an Agent Server (see Fig. 7), with concrete sponsors including a containerized server (`DockerWorkspace`) or a API-managed runtime (`APIRemoteWorkspace`). The factory `Workspace(...)` resolves to local when only `working_dir` is provided and to remote when `host`/runtime parameters are present, ensuring the agent code remains unchanged across environments.

## 5 Reliability and Evaluation

We assess the reliability and performance of the OpenHands Agent SDK through two complementary processes: continuous testing and benchmark evaluation. The continuous testing pipeline combines both programmatic and LLM-based tests and runs automatically on every pull request and once per day. It checks that the SDK behaves consistently across multiple language models, catching regressions in reasoning, tool use, and state management early. These automated tests cost only $0.5–$3 per full run and complete in less than 5 minutes. In parallel, benchmark evaluations measure the SDK's overall capabilities on standardized agentic tasks, providing a broader view of model quality and system performance.

```python
class BaseWorkspace(ABC):
    def __enter__(self) -> "BaseWorkspace"
    def __exit__(self, exc_type: Any, exc_val: Any, exc_tb: Any) -> None
    def execute_command(self, command: str) -> CommandOutput:
        """Run shell command in workspace environment."""
    def file_upload(self, path: str, content: bytes) -> None:
        """Write file to workspace filesystem."""
    def file_download(self, path: str) -> bytes:
        """Read file from workspace filesystem."""

class LocalWorkspace(BaseWorkspace):
    def execute_command(self, cmd: str) -> CommandOutput:
        return subprocess.run(cmd, shell=True, capture_output=True)

class RemoteWorkspace(BaseWorkspace):
    def execute_command(self, cmd: str) -> CommandOutput:
        return self._api_client.post("/execute", json={"command": cmd})

# Workspace factor class
class Workspace:
    def __new__(cls, *, host=None, working_dir="workspace/project", api_key=None):
        if host is provided:
            return RemoteWorkspace(working_dir, host, api_key)
        else:
            return LocalWorkspace(working_dir)
```

**Figure 7** Workspace interface. Implementations handle environment details.

```python
class BaseIntegrationTest(ABC):
    def setup(self): """Prepare test environment"""
    def tools(self) -> list[Tool]: """Agent capabilities"""

    def verify_result(self) -> bool: """Check success"""
    def run(self, instruction: str) -> bool:
        self.setup()
        conversation = Conversation(Agent(llm, self.tools), workspace)
        conversation.send_message(instruction)
        conversation.run()
        return self.verify_result()
```

**Figure 8** Integration test framework. Subclasses implement `setup()`, `tools`, and `verify_result()` for specific scenarios.

## 5.1 Continuous Quality Assurance

The SDK employs a three-tier testing strategy that balances coverage, cost, and depth:

- **Programmatic Tests** — Run on every commit. These tests mock llm calls and verify core logic, data flow, and API contracts within seconds. Mocking allows quicker feedback, ensuring that most regressions are caught before any external API calls are made.
- **LLM-based Tests** — Include both integration and example tests (see below). Executed daily and on-demand for pull requests. These tests use real models (Claude Sonnet 4.5, GPT-5 Mini, DeepSeek Chat) to validate reasoning, tool invocation, and environment stability. Each run costs $0.5–$3 and completes in less than 5 minutes.
- **Benchmark Evaluation** — On-demand, high-cost evaluations ($100–1000, hours per run) that measure comprehensive agent capabilities on academic datasets.

**Integration tests** cover multiple scenario-based workflows (e.g., file manipulation, command execution, git operations, and browsing), while **example tests** periodically run all SDK examples (custom tools, MCP integration, persistence, async execution, routing, etc) to ensure end-to-end reliability. The suite is continuously expanded as new agent behaviors and failure patterns are discovered, improving coverage and regression sensitivity over time. On-demand execution for these LLM-based tests further optimizes CI/CD cost: integration tests target high-risk changes, example tests cover user-facing modules, and daily runs track regressions across overall codebase updates.

## 5.2  Benchmarks

The SDK offers built-in support for a diverse range of academic benchmarks evaluating agentic capability.

| Benchmark | Model | Performance |
|---|---|---|
| SWE-Bench (Jimenez et al., 2024) | Claude Sonnet 4.5 | 72.8% |
| | Claude Sonnet 4 | 68.0% |
| | GPT-5 (reasoning=high) | 68.8% |
| | Qwen3 Coder 480B A35B | 65.2% |
| GAIA (Mialon et al., 2023a) | Claude Sonnet 4 | 57.6% |
| | Claude Sonnet 4.5 | 67.9% |
| | GPT-5 (reasoning=high) | 62.4% |
| | Qwen3 Coder 480B A35B | 41.2% |

**Table 2** Benchmark performance of OpenHands Software Agent SDK. Evaluations were performed at commit `54c5858` of the SDK and commit `88f1d80` of the benchmarks.

As shown in Tab. 2, the SDK demonstrates competitive performance on software engineering and general agent benchmarks.

On SWE-Bench Verified (Jimenez et al., 2024) that measures an agent's ability in software engineering tasks, SDK achieves 72% resolution rate using Claude Sonnet 4.5 with extended thinking; On the GAIA (val set, Mialon et al. (2023b)) that measuresan agent's generic computer task-solving capability, SDK achieves 67.9% accuracy with Claude Sonnet 4.5, demonstrating effective multi-step reasoning and tool use. Additionally, Qwen3 Coder 480B (Yang et al., 2025), a strong open model for coding, also stays at a competitive score of 41.21%. These results are slightly better than those of OpenHands-Versa (Soni et al., 2025), validating that the SDK's architecture does not sacrifice agentic capability and maintains performance competitive with research-focused systems.

## 6  Related Work

Recent agent SDKs share common goals of tool orchestration, state management, and production deployment. General-purpose frameworks such as LangChain and LangGraph focus on compositional pipelines and stateful graph execution with durable checkpoints for long-running reasoning workflows (LangChain, 2025; Team, 2025). Provider SDKs from OpenAI and Anthropic emphasize production orchestration, guardrails, and handoffs within their own model ecosystems (OpenAI, 2025; Anthropic, 2025a), while Google's Agent Development Kit (ADK) targets model-agnostic agents integrated with Vertex AI for managed deployment (Google, 2025).

In contrast, the *OpenHands Software Agent SDK* is a fully open-source, vendor-agnostic platform purpose-built for *software engineering agents*. It couples an event-sourced state model and immutable configuration with sandboxed execution, a built-in REST/WebSocket server, and workspace-level remote interfaces (VS Code, VNC, browser). These design choices enable reproducible, deterministic execution and seamless transition from local prototyping to production-scale deployments—capabilities not addressed by existing general-purpose SDKs. A feature-level comparison with major agent SDKs is provided in Tab. 3.

## 7  Conclusion

The OpenHands Software Agent SDK bridges the gap between rapid local prototyping and production deployment for software engineering agents through a stateless, event-sourced, and composable architecture spanning four packages (SDK, Tools, Workspace, Server). From our experience scaling OpenHands from research prototype to production system, we learned that strict separation between core agent logic and

| | OpenAI Agents SDK | Claude Agent SDK | Google Agent Development Kit | OpenHands Software Agent SDK |
|---|---|---|---|---|
| *Standalone SDK Features* | | | | |
| MCP Support | ✓ | ✓ | ✓ | ✓ |
| Custom Tools | ✓ | ✓ | ✓ | ✓ |
| History Persistence & Restore | ✓ | ✓ | ✓ | ✓ |
| Sub-agent Delegation | ✓ | ✓ | ✓ | ✓ |
| Model Agnostic (100+ LLMs) | ✓ | ✗ | ✓ | ✓ |
| Multi-LLM Routing | ✗ | ✗ | ✗ | ✓ |
| Conversation Cost & Token Tracking | ✓ | ✓ | ✗ | ✓ |
| Pause/Resume Agent Execution | ✓ | ✗ | ✗ | ✓ |
| Native Support for Non-function-calling models | ✗ | ✗ | ✗ | ✓ |
| Security Analyzer for Agent Action | ✗ | ✗ | ✓ | ✓ |
| Action Confirmation Policies | ∼ | ✗ | ∼ | ✓ |
| Context Files (e.g., repo.md, AGENTS.md) | ✗ | ✓ | ✗ | ✓ |
| Agent Skills | ✗ | ✓ | ✗ | ✓ |
| Context Condensation | ✗ | ✓ | ✗ | ✓ |
| TODO List Planner | ✗ | ✓ | ✗ | ✓ |
| Tmux-based Interactive Bash Terminal | ✗ | ✗ | ✗ | ✓ |
| Auto-Generated Conversation Titles | ✗ | ✗ | ✗ | ✓ |
| Secrets Management with Auto-Masking | ✗ | ✗ | ✗ | ✓ |
| Agent Stuck Detection | ✗ | ✗ | ✗ | ✓ |
| Long-term Memory across Sessions | ✗ | ✗ | ✓ | ✗ |
| *Production Server Features* | | | | |
| Builtin REST+WebSocket Server | ✗ | ✗ | ✗ | ✓ |
| Session-Based Authentication | ✗ | ✗ | ✗ | ✓ |
| Builtin Remote Agent Execution | ✗ | ✗ | ✗ | ✓ |
| Agent Environment Sandboxing | ✗ | ✗ | ∼ | ✓ |
| VNC Desktop for Agent Workspace | ✗ | ✗ | ✗ | ✓ |
| VSCode Web for Agent Workspace | ✗ | ✗ | ✗ | ✓ |
| Builtin Chromium Browser for Agent | ✗ | ✗ | ✗ | ✓ |
| *Continous Quality Assurance* | | | | |
| Programmatic tests | ✓ | ✓ | ✓ | ✓ |
| Strong Type Checking | ✓ | ✓ | ✓ | ✓ |
| Frequent LLM-based tests | ✗ | ✓ | ✓ | ✓ |
| Built-in Academic Benchmark Evaluation | ✗ | ✗ | ✗ | ✓ |

**Table 3** Agent SDK features comparisons (✓= full support, ∼= partial, ✗= absent; as of Oct 29, 2025). OpenHands software agent SDK uniquely combines: **(i)** native remote execution with secure sandboxing (vs. external orchestration for OpenAI and local-only for Claude), **(ii)** a built-in production server with REST+WebSocket APIs and workspace access (vs. library-only SDKs), **(iii)** LLM-powered action-level security analysis (vs. manual or external guardrails), and **(iv)** model-agnostic multi-LLM routing with first-class support for non-function-calling models.

applications is essential for maintainability; event-sourced state enables reproducibility and fault recovery; immutable component design prevents configuration drift; and a unified execution model that supports both local development and sandboxed deployment streamlines the transition from experimentation to production. These lessons collectively shaped a system that supports reliable, extensible operation across heterogeneous environments. Evaluation on SWE-Bench Verified and GAIA benchmarks confirms strong performance and consistency across diverse model backends, validating the SDK as a robust foundation for both research and industrial-scale deployment.

# References

Anthropic. Claude agent sdk: Overview and python sdk. https://anthropic.mintlify.app/en/api/agent-sdk/overview, 2025a. Accessed: 2025-10-29.

Anthropic. Claude code: An agentic coding tool that lives in your terminal, 2025b. URL https://github.com/anthropics/claude-code. GitHub repository, accessed 2025-10-28.

BerriAI. Litellm: Call 100+ llm apis in openai format. https://github.com/BerriAI/litellm, 2024. Accessed: 2025-01-06.

Cognition AI. Devin: Ai software engineer. https://www.cognition.ai/devin, 2024. Accessed: 2025-01-06.

Cursor Team. Cursor: The ai-first code editor. https://www.cursor.com, 2024. Accessed: 2025-01-06.

GitHub. Github copilot: Your ai pair programmer. https://github.com/features/copilot, 2021. Accessed: 2025-01-06.

Google. Agent development kit (adk). https://google.github.io/adk-docs/, 2025. Accessed: 2025-10-29.

Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. SWE-bench: Can language models resolve real-world github issues? In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=VTF8yNQM66.

LangChain. Langchain: Runnables and the langchain expression language (lcel). https://api.python.langchain.com/en/latest/core/runnables.html, 2025. Accessed: 2025-10-29.

Guangya Liu. Agents.md: The readme for your ai coding agents. https://research.aimultiple.com/agents-md/, August 2025. Accessed: 2025-01-06.

MCP Team. Model context protocol (mcp)? https://modelcontextprotocol.io, 2025. Accessed: 2025-10-02.

Grégoire Mialon, Clémentine Fourrier, Thomas Wolf, Yann LeCun, and Thomas Scialom. Gaia: a benchmark for general ai assistants. In *The Twelfth International Conference on Learning Representations*, 2023a.

Grégoire Mialon, Clémentine Fourrier, Craig Swift, Thomas Wolf, Yann LeCun, and Thomas Scialom. Gaia: a benchmark for general ai assistants, 2023b. URL https://arxiv.org/abs/2311.12983.

Graham Neubig. One year of openhands: A journey of open source ai development. *All Hands AI Blog*, March 2025. URL https://www.all-hands.dev/blog/one-year-of-openhands-a-journey-of-open-source-ai-development.

OpenAI. Openai agents sdk. https://github.com/openai/openai-agents-python, 2024. Accessed: 2025-01-06.

OpenAI. Agents sdk and guide. https://platform.openai.com/docs/guides/agents, 2025. Accessed: 2025-10-29.

Pydantic Team. Unions — discriminated unions. https://docs.pydantic.dev/latest/concepts/unions/#discriminated-unions, 2025. Accessed: 2025-10-29.

Stuart Russell and Peter Norvig. *Artificial Intelligence: A Modern Approach*. Prentice Hall Press, USA, 3rd edition, 2009. ISBN 0136042597.

Calvin Smith. Openhands context condensensation for more efficient ai agents. *All Hands AI Blog*, April 2025. URL https://openhands.dev/blog/openhands-context-condensensation-for-more-efficient-ai-agents.

Aditya Bharat Soni, Boxuan Li, Xingyao Wang, Valerie Chen, and Graham Neubig. Coding agents with multimodal browsing are generalist problem solvers, 2025. URL https://arxiv.org/abs/2506.03011.

LangGraph Team. Langgraph documentation (durable execution, deployment, server/cloud). https://docs.langchain.com/oss/python/langgraph/, 2025. Accessed: 2025-10-29.

Lei Wang, Chengbang Ma, Xueyang Feng, Zeyu Zhang, Hao ran Yang, Jingsen Zhang, Zhi-Yang Chen, Jiakai Tang, Xu Chen, Yankai Lin, Wayne Xin Zhao, Zhewei Wei, and Ji rong Wen. A survey on large language model based autonomous agents. *ArXiv*, abs/2308.11432, 2023. URL https://api.semanticscholar.org/CorpusID:261064713.

Xingyao Wang, Boxuan Li, Yufan Song, Frank F. Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan, Yueqi Song, Bowen Li,

Jaskirat Singh, Hoang H. Tran, Fuqiang Li, Ren Ma, Mingzhang Zheng, Bill Qian, Yanjun Shao, Niklas Muennighoff, Yizhe Zhang, Binyuan Hui, Junyang Lin, Robert Brennan, Hao Peng, Heng Ji, and Graham Neubig. Openhands: An open platform for AI software developers as generalist agents. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=OJd3ayDDoF.

An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report, 2025. URL https://arxiv.org/abs/2505.09388.

# Appendix

## A    SDK Feature Comparison Details

This appendix provides documentation sources and implementation details for capabilities listed in Tab. 3.

### A.1    Model Context Protocol (MCP) Support

All four SDKs provide support for the Model Context Protocol (MCP Team, 2025), which standardizes how agents interact with external tools and data sources:

**OpenHands:** Full MCP integration with OAuth, supporting both local and remote MCP servers. (`https://docs.openhands.dev/sdk/guides/mcp`)

**OpenAI:** Native MCP support documented at `https://openai.github.io/openai-agents-python/mcp/`

**Claude:** MCP integration guide at `https://docs.claude.com/en/api/agent-sdk/mcp`

**Google:** MCP usage instructions at `https://google.github.io/adk-docs/mcp/`

### A.2    Model Agnostic Support (100+ LLMs)

**OpenHands Agent SDK:** Native support for 100+ LLMs via LiteLLM integration (BerriAI, 2024), enabling seamless switching between providers (OpenAI, Anthropic, Google, local models) without code changes.

**Claude Agent SDK:** Locked to Anthropic's Claude models only. The SDK architecture is tightly coupled to Anthropic's API and does not support model-agnostic execution (`https://docs.anthropic.com/en/docs/build-with-claude/agent-sdks`).

**OpenAI Agents SDK & Google ADK:** Support multiple models via LiteLLM extension. While the SDKs are designed primarily for OpenAI and Gemini models respectively, they can be extended to support other providers through community contributions (`https://github.com/openai/openai-agents-python`, `https://google.github.io/adk-docs/agents/models/`).

### A.3    Multi-LLM Routing Support

**OpenHands Agent SDK:** Provides built-in LLM routing with `MultimodalRouter` and `RandomRouter` classes (`https://github.com/OpenHands/software-agent-sdk/tree/main/openhands-sdk/openhands/sdk/llm/router`). Features include:

- Automatic routing between models based on content type (text vs. images)
- Cost-based routing to optimize expenditure
- Transparent fallback mechanisms

Example usage in `examples/19_llm_routing.py` demonstrates routing between Claude Sonnet (primary) and Mistral (secondary) based on content type.

**OpenAI, Claude Agent SDK & Google ADK:** No built-in routing capability. Developers must manually implement model selection logic at the application layer.

### A.4    Pause/Resume Agent Execution

**OpenHands Agent SDK:** Built-in pause/resume mechanism with `AgentExecutionStatus` tracking (`https://github.com/OpenHands/software-agent-sdk/blob/main/openhands-sdk/openhands/sdk/conversation/state.py`). Features include:

- `conversation.pause()`: Gracefully interrupt agent execution
- `conversation.resume()`: Continue from last checkpoint
- Status tracking: RUNNING, PAUSED, FINISHED, ERROR
- Example: `examples/09_pause_example.py` demonstrates `KeyboardInterrupt` handling

This capability is critical for interactive applications where users may need to interrupt long-running agents, provide additional context, or change direction mid-execution.

**OpenAI, Claude Agent SDK & Google ADK:** No built-in pause/resume mechanism. Developers must implement custom state management and checkpointing.

## A.5    Support for Non-function-calling models

**OpenHands Agent SDK:** Implements automatic fallback to prompt-based function calling for models that don't natively support tool calling ([https://github.com/OpenHands/software-agent-sdk/blob/main/openhands-sdk/openhands/sdk/llm/mixins/non_native_fc.py](https://github.com/OpenHands/software-agent-sdk/blob/main/openhands-sdk/openhands/sdk/llm/mixins/non_native_fc.py)). This enables transparent tool use across 100+ LLMs.

**OpenAI, Claude Agent SDK & Google ADK:** Only support models with native function calling. No fallback mechanism for models lacking this capability.

## A.6    Security Analyzer for Agent Action & Action Confirmation Policies

**OpenHands Agent SDK:** Implements multi-layer security where the LLM analyzes each action before execution, assigning risk levels (LOW, MEDIUM, HIGH) based on system impact:

- **Risk assessment:** The LLM is prompted to include a security risk field when generating tool calls, evaluating each proposed action (bash command, file operation, API call) by reasoning about its effects: "This command installs packages with `sudo`—HIGH risk because it modifies system state."
- **Confirmation policies:** Based on risk level, the system triggers appropriate workflows, such as:
  - LOW: Execute immediately
  - HIGH: Pause execution, request human confirmation
- **Customizable policies:** Developers define `ConfirmationPolicy` classes to match their risk tolerance (e.g., auto-approve in sandbox environment, require confirmation in production).
- **Context-aware:** The LLM has access to conversation history and can reason about cumulative risk (e.g., "This is the third `rm -rf` command in this conversation—escalating risk level").

This approach balances safety with usability: routine operations proceed smoothly, while dangerous actions trigger safeguards. It's particularly valuable for production deployments where agents handle sensitive data or have elevated privileges.

**Google ADK:** Even through the ADK doesn't provide a direct implementation, primitives like plugins and callbacks allow users to easily create LLM-based safety guardrails, in which models analyze user inputs as well as tool input and output. ([https://google.github.io/adk-docs/safety/#callbacks-and-plugins-for-security-guardrails](https://google.github.io/adk-docs/safety/#callbacks-and-plugins-for-security-guardrails))

**OpenAI & Claude Agent SDKs:** Neither provides automated security analysis. Developers must implement manual guardrails by:

- Writing custom validation logic for each tool
- Hardcoding blacklists of dangerous commands (brittle, easily bypassed)
- Building separate confirmation workflows (complex, error-prone)

This places the entire security burden on application developers, who may lack expertise in adversarial scenarios or prompt injection attacks.

## A.7  Static Context Files

**OpenHands:** Supports AGENTS.md, repo.md files and other third party files that provide repository-wide guidelines, coding standards, etc.

**Claude:** Supports CLAUDE.md files for project-specific context (https://docs.claude.com/en/api/agent-sdk/overview).

**OpenAI Agents SDK & Google ADK:** No support for static context files. Note: AGENTS.md is a feature of OpenAI Codex (the IDE integration), not the Agents SDK (Liu, 2025).

## A.8  Context Condensation

**OpenHands:** Automatically condenses conversation history when context windows approach limits, using summarization to preserve critical information while reducing token count.

**Claude:** Offers automatic, *server-side* context editing with configurable strategies (particularly the clear_tool_uses_20250919 and clear_thinking_20251015 strategies) that remove stale tool results and thinking blocks based on token thresholds, allowing the conversation to continue without hitting context limits.

**OpenAI Agents SDK & Google ADK:** No automatic context condensing feature; developers must manually manage context window limits.

## A.9  Auto-Generated Conversation Titles

**OpenHands Agent SDK:** Automatic conversation title generation using LLM-based summarization of the first user message (https://github.com/OpenHands/software-agent-sdk/blob/main/openhands-sdk/openhands/sdk/conversation/title_utils.py). Features include:

- Category detection with emojis (e.g., 🐛 bugfix, ⭐ features, 📄 documentation)
- Concise, descriptive titles for conversation history UIs
- Configurable max length to prevent overflow in UI elements

**OpenAI, Claude Agent SDKs & Google ADK:** No automatic title generation. Developers must implement custom title extraction or use first message text verbatim.

## A.10  Secrets Management with Auto-Masking

**OpenHands Agent SDK:** Automatic secrets management system (https://github.com/OpenHands/software-agent-sdk/blob/main/openhands-sdk/openhands/sdk/conversation/secrets_manager.py) that:

- Detects secrets in bash commands (API keys, passwords, tokens)
- Automatically injects secrets as environment variables in sandboxed execution
- Masks secret values in logs and LLM context to prevent leakage
- Supports pattern-based detection and explicit secret registration

This prevents accidental exposure of sensitive credentials in conversation history, logs, or debugging outputs.

**OpenAI, Claude Agent SDKs & Google ADK:** No built-in secrets management. Developers must manually implement masking and secure credential handling.

## A.11  Stuck Detection

**OpenHands:** Implements automatic detection of pathological agent states, including:

- Infinite loops (same action repeated without progress, same agent errors)

- Redundant tool calls (querying the same information repeatedly)

When stuck behavior is detected, the system can automatically intervene by terminating execution to prevent resource waste.

**OpenAI, Claude Agent SDKs & Google ADK:** No stuck detection mentioned in documentation; developers must implement custom monitoring.

## A.12  Long-term Memory Across Session

**Google ADK:** Provides a `MemoryService` interface with two implementations: in-memory for prototyping and production use with managed persistence via Vertex AI API. The service allows agents to ingest completed sessions into long-term storage and retrieve relevant information from past conversations through search queries, typically accessed via tool calls. (https://google.github.io/adk-docs/sessions/memory/)

## A.13  Built-in REST + WebSocket Server

**OpenHands:** Includes a production-ready FastAPI server with both REST endpoints (for request/response operations) and WebSocket support (for real-time streaming). This eliminates the need for developers to build custom server infrastructure.

**OpenAI, Claude Agent SDK & Google ADK:** Framework-only, no built-in server; developers must build their own server layer.

## A.14  Remote Execution Architecture

**OpenHands Agent SDK:** Provides native support for distributed deployments where agents run on client machines while tools execute in remote sandboxed environments. This architecture enables:

- **Separation of concerns:** Agent logic runs locally (low latency, private), while tool execution runs remotely (isolated, scalable)
- **Resource flexibility:** High-compute tools (e.g., code compilation, browser automation) run on remote servers without client hardware constraints
- **Security isolation:** Untrusted code executes in remote containers, protecting client machines
- **Multi-tenancy:** Multiple users share remote infrastructure while maintaining workspace isolation

The implementation uses WebSocket connections for bidirectional communication between agent and remote runtime, with automatic reconnection and state synchronization. Configuration is simple: specify `remote_-runtime_url` and the SDK handles the rest.

**OpenAI Agents SDK:** Requires external Temporal workflow orchestration for production deployment (OpenAI, 2024). This adds significant infrastructure complexity:

- Developers must set up and maintain Temporal servers
- Workflow definitions require learning Temporal's API and concepts
- Integration overhead: agents, Temporal workers, and application servers must be coordinated
- No built-in solution for simple remote execution scenarios

While Temporal provides powerful orchestration capabilities (retries, workflows, versioning), this is overkill for many use cases and creates a steep learning curve.

**Claude Agent SDK & Google ADK:** Supports local execution only. Agents and tools must run on the same machine where the SDK is invoked. This limits:

- Scalability: Cannot offload computation to remote servers
- Security: All tool execution happens on the client machine

• Deployment flexibility: Cannot easily create multi-tenant services

## A.15 Agent Environment Sandboxing

**OpenHands:** First-class Docker runtime support that isolates code execution in containers, providing:

• Complete filesystem isolation
• Network policy enforcement
• Resource limits (CPU, memory, disk)
• Easy reproducibility across environments

**Google ADK**: Cloud-dependent sandboxing via Agent Engine with process-level isolation and persistent state, but requires Google Cloud project setup. ([https://google.github.io/adk-docs/tools/google-cloud/code-exec-agent-engine/](https://google.github.io/adk-docs/tools/google-cloud/code-exec-agent-engine/))

**OpenAI & Claude Agent SDKs:** No sandboxing capability; all execution occurs in the local environment, requiring manual security measures.

## A.16 Interactive Workspace Access

**OpenHands Agent SDK:** The production server exposes the agent's sandbox through three mechanisms, enabling real-time monitoring and intervention during agent execution:

1. **VNC Desktop:** Full desktop GUI access via VNC protocol. Users can observe the agent's file system, terminal sessions, and running processes in real time. Useful for debugging and understanding agent behavior on a higher level.
2. **Browser Access:** Built-in Chromium browser accessible via the web UI in non-headless mode. Users can see what the agent sees when browsing web pages, enabling debugging of web scraping or navigation tasks.
3. **VSCode Web:** Full VSCode editor embedded in the workspace. Users can:
   • Edit files the agent is working with
   • Review code changes in real time
   • Intervene manually when the agent gets stuck
   • Collaborate with the agent

These capabilities are essential for production debugging and human-in-the-loop workflows. They transform agents from black boxes into observable, interactive systems.

**OpenAI, Claude Agent SDK & Google ADK:** None of these provides built-in workspace interaction capabilities. Agents are opaque processes—once started, there's no way to observe their workspace, inspect files, or intervene without stopping execution and examining logs. This makes debugging production failures significantly harder.

## A.17 Quality Assurance Practices

The comparison table reveals significant differences in how each SDK approaches testing and quality assurance, particularly regarding the use of real LLMs versus mocked responses in their test suites.

**OpenHands Agent SDK** employs a unique dual-track testing strategy that sets it apart from other SDKs:

• **Unit tests with mocked LLMs:** Fast, deterministic tests for core logic, tool functions, and state management. These run on every commit without incurring API costs.
• **Integration tests with real LLMs:** The SDK runs frequent integration tests against actual language models to verify agent behavior, prompt effectiveness, and tool orchestration. These tests are strategically scheduled (e.g., nightly builds, pre-release validation) to balance coverage with cost.

- **Academic benchmark integration:** OpenHands is the only SDK with built-in evaluation against standardized benchmarks (Jimenez et al., 2024), enabling quantitative performance tracking across releases. This provides empirical evidence of agent capabilities beyond anecdotal testing.

**OpenAI Agents SDK** follows conventional software testing practices with comprehensive unit tests and strong typing but notably lacks frequent LLM-based integration testing. Based on our analysis of their repository[1]:

- Tests primarily use mocked responses to avoid API costs
- Integration tests, when present, are limited and typically run manually
- No systematic evaluation against academic benchmarks
- Relies on community feedback for real-world validation

**Claude Agent SDK** presents an unusual testing approach due to its architecture (Python SDK → subprocess → Claude Code CLI → API):

- Unit tests exist but are complicated by the subprocess communication layer
- The CLI dependency makes mocking challenging, often requiring real API calls even for basic tests
- Platform-specific issues (Windows vs. Linux) have plagued the test suite[2]
- No benchmark evaluation framework

**Google ADK** adopts a hybrid approach with its evaluation-centric framework:

- Mock tools for deterministic component testing
- Real Gemini API calls for agent logic validation
- Structured "evalsets" for systematic performance assessment[3]
- Focus on trajectory evaluation (comparing expected vs. actual tool usage sequences)

---

[1] https://github.com/openai/openai-agents-python
[2] See issues #208 and #39 in https://github.com/anthropics/claude-agent-sdk-python
[3] Documentation at https://google.github.io/adk-docs/evaluation/