**⊛ ChatGPT**

# Application Development with Autonomous AI Agents

## Introduction

AI-driven software development has evolved from prompting single models to orchestrating **multiple specialized AI agents** collaborating on complex tasks. By 2025, success hinges on managing these agent teams effectively rather than relying on one monolithic assistant [1] . Frameworks like **LangGraph** enable stateful, fault-tolerant *multi-agent workflows* with memory and retries – a "Terraform for agent orchestration," as enterprises seek AI-native systems with robust automation [2] [3] . In practice, well-orchestrated agents can parallelize work, cross-verify each other's outputs, and leverage diverse skills to achieve results that single models struggle to deliver [4] [5] . This report presents a **methodology** for agent-assisted application development, covering the full lifecycle from planning to deployment. It details **agent roles** and communication protocols, compares leading **orchestration frameworks**, and outlines best practices in **context management**, **structured I/O contracts**, and **safe, deterministic execution**. We also provide guidance on modular, testable build cycles with validation gates and recovery mechanisms, integration tips for tools like GitHub and VS Code, and templates for schemas, role taxonomy, architecture patterns, observability, and retrospectives.

## Lifecycle of Agent-Assisted Development

Building an application with AI agents involves iterative phases, each with targeted agent support:

- **Planning & Design:** A *Strategist/Architect agent* analyzes requirements and devises a high-level solution approach. It breaks the project into phases and tasks, accounting for safety and verification steps. The plan should include testing checkpoints and continuity notes so that any agent joining later can get up to speed [6] [7] . The strategist focuses on *clarity* and *risk mitigation*, producing a roadmap with well-defined milestones (e.g. feature modules, integration steps) and "hold points" for validation.

- **Task Decomposition:** Once the roadmap is set, the strategist (or an Orchestrator) generates an **ultra-atomic task list** for the next phase [8] . Tasks are kept **small and self-contained** – each fitting within an LLM's context window and yielding a tangible output (code for a function, test results, etc.). Breaking work into bite-sized chunks prevents context overflow and reduces the chance of agent confusion or hallucination [9] [10] . After drafting the list, a validation step can have the agent double-check for any gaps, unnecessary complexity, or potential failure points [11] , adjusting tasks until there's high confidence (often aiming for >95% estimated success) [12] .

- **Implementation (Coding):** A specialized *Execution agent* (Coder) tackles the tasks in sequence (or in parallel if independent). It writes code or content according to the specifications. Crucially, the execution is **guarded by incremental checkpoints**: the agent should stop after completing a logical

sub-task and report status [13] . This allows the orchestrator or other agents to review progress before moving on. During coding, the agent adheres to the I/O contracts (e.g. function signatures, file formats) defined in the plan. Some workflows even instruct the coder agent to internally self-review or use a rubric, ensuring it meets quality criteria before "submitting" code [14] . By treating each step as a transaction that must be verified, the process avoids compounding errors.

· **Quality Control & Validation:** After the execution agent produces an output, one or more agents perform QA. A *Reviewer agent* (QC) conducts a **code review or analysis** of the changes, looking for bugs, inconsistencies, or unmet requirements [15] . It may generate a list of issues ranked by severity without directly fixing them [16] [17] . Next, a *Tester agent* runs **module-level tests** [18] . This can include running automated tests, validating that functions return expected results, or checking that the overall module behaves as intended. The tester should log outcomes (perhaps in a test report file or console output) for the orchestrator to evaluate. Finally, a *Validation agent* performs an **independent verification**: it cross-checks that the implementation actually matches what the strategist planned and what the execution agent claimed to do [19] . This agent essentially asks "Did we build what we said we would, and is it correct?" – verifying all fixes were applied and that the system meets the acceptance criteria with high confidence [20] . Only after these checks pass does the workflow proceed. If issues are found, the system can loop back: e.g. the orchestrator might assign the execution agent to fix specific bugs, or even engage a *different* agent (perhaps a specialist) to assist, before re-running tests. These **validation gates** ensure that errors are caught early, preventing faulty code from propagating to later stages [9] . Notably, multi-agent setups inherently provide *built-in quality checks* – different agents bring different perspectives and can catch each other's mistakes [21] .

· **Deployment & Documentation:** With features implemented and validated, the focus shifts to integrating the work into the broader project and keeping records. A *Documentation agent* (Scribe) updates **continuity documentation** – essentially the project logs, README, or design docs – to reflect the current status, decisions made, and next steps [7] . This is optimized for machine readability (structured and concise), since future agents or sessions will rely on it to understand context [22] . The documentation agent might also compile **lessons learned** during this iteration (e.g. challenges faced, how they were overcome, suggestions for improvement) [23] – providing data for process refinement. Additionally, code comments or usage docs can be generated for any new modules (ensuring maintainability). Finally, changes are **deployed or saved**: the documentation or orchestrator agent commits code to a repository and may push to GitHub or another version control system [24] . Regular check-ins not only serve as backups against crashes but also enable collaboration and auditability. In more advanced pipelines, an agent could trigger a CI/CD workflow or send a deployment request after passing all tests. At minimum, the orchestrator confirms that the new code is integrated into the main codebase safely.

· **Monitoring & Iteration:** (Beyond initial deployment) Agents can assist in monitoring application performance or user feedback post-deployment, though this strays into operations. For completeness, one might imagine an *Observability agent* watching logs or a *Feedback analyst* agent summarizing user reviews for the team. These are emerging areas; in practice many teams handle ops externally, but the framework is extensible. Crucially, the development process repeats for each feature or improvement in an **iterative cycle**, using prior lessons to refine strategies. A retrospective step (often led by the human team or a designated agent) evaluates what to adjust in the next sprint

– e.g. maybe the task breakdown could be improved or a new test agent is needed. We provide a checklist for such evaluation in the Appendix.

Throughout this lifecycle, an **Orchestrator agent** (or a simple central controller) oversees transitions between phases. It ensures that each stage's exit criteria are met before moving to the next. For example, the orchestrator would only proceed from coding to testing after the strategist's plan and task list received a thumbs-up (high confidence) and would only deploy after all tests pass and documentation is updated. If any step fails, the orchestrator can re-route the workflow (e.g. prompt the strategist to revise the plan or have the coder attempt a fix). In Microsoft's *Magentic-One* system, for instance, a lead Orchestrator agent maintains a **task ledger** and a *progress ledger*, tracking what's done and self-reflecting on completeness; if agents get stuck or errors occur, the Orchestrator adapts the plan or assigns new tasks to recover [25] [26]. This kind of supervised coordination is key to robust multi-agent development.

## Agent Roles and Responsibilities

Proper **role decomposition** is the backbone of effective multi-agent systems [9]. Each agent is assigned a clear role with defined responsibilities, success criteria, and failure modes. Below we outline a common taxonomy of roles in an AI development team and how they interact, followed by a more detailed table in the **Agent Taxonomy** document.

- **Orchestrator / Supervisor** – **Responsibility:** High-level coordination of all other agents and overall workflow control. The orchestrator parses the user's request or project goal and decides which agents to invoke, in what order [27]. It manages the "conversation" among agents, making sure outputs from one flow into the next appropriately. It also enforces validation gates: e.g. verifying that a plan exists before coding, or that tests passed before deployment. **Success metrics:** The orchestrator succeeds if the multi-agent system completes tasks efficiently and correctly – i.e., tasks are assigned to the right agents, bottlenecks are minimized, and errors are caught and handled. One simple metric is *cycle efficiency* (how many loops/retries were needed) and *goal completion rate*. **Failure modes:** If the orchestrator misroutes tasks or fails to detect an agent's failure, the workflow can stall or produce faulty output. For instance, not handling a non-responsive agent (infinite loop) is an orchestrator failure. Also, overly rigid or incorrect planning by the orchestrator can lead the team astray (garbage in, garbage out). Robust orchestrators use strategies like hierarchical planning and can fall back or ask for help when needed [28] [25]. Often, the orchestrator role may be partially fulfilled by a human overseer or a framework's runtime logic rather than an LLM.

- **Strategist / Architect** – **Responsibility:** Interpreting objectives and requirements, then formulating a solution approach. This agent conducts "big picture" thinking: breaking the project into components, choosing architectural patterns, and drafting implementation plans. It also defines communication protocols (e.g. what each agent's input/output should look like) and ensures continuity information is recorded. The strategist might produce design docs, flowcharts, or pseudocode outlining how the problem will be solved in stages [29] [7]. **Success metrics:** A good plan covers all key requirements with clear, feasible steps and includes testing/validation embeds. One can measure plan quality by subsequent smooth execution – e.g., low incidence of major revisions needed, and high confidence ratings from the validation agent or orchestrator. The strategist's plan should also minimize risk (no obvious gaps or overly large tasks). **Failure modes:** Poor planning leads to scope creep, missing features, or tasks that exceed the context window or agent capabilities. If the strategist under-specifies tasks, the execution agent might misinterpret

them; if it over-specifies or includes too much in one step, the agent may run out of context or make mistakes. Another failure mode is not anticipating dependency or integration issues (e.g., two agents working on incompatible assumptions). To mitigate this, strategists often include **checkpoints and fallback plans** (for example, "if module X fails test, do Y"). A strategist agent might loop on its own output – reviewing and refining the plan – until a confidence threshold (like ">95% confident in plan viability") is reached [30] [12] .

- **Execution / Builder** – **Responsibility:** Constructing the actual application artifacts (code, configuration, content) according to the plan. This agent is essentially the "developer," writing code functions, creating files, modifying data, etc. It may use tools (compilers, linters, API calls) as needed to fulfill tasks. In advanced frameworks, the execution agent can spawn code in a sandbox or call lower-level functions. **Success metrics:** Measurable outputs like *code correctness* (does it run and meet specs), and *efficiency* (did it complete within resource limits). The execution agent is successful when it implements each task without introducing bugs and while following the specified interfaces/ contracts. Ideally, its output should pass all tests on the first run – though in practice some iterations are normal. **Failure modes:** The execution agent can fail by producing **syntax errors**, *logical errors*, or outputs that violate the expected schema or format. To guard against this, many systems enforce structured output for code actions (e.g., having the agent output a JSON or function call representing the code change, which is then applied by a deterministic interpreter) [31] [32] . If the agent does free-form generation, it might hallucinate functions or use incorrect APIs. Another failure mode is **overstepping scope** – e.g. modifying parts of the project it wasn't supposed to (hence the importance of context and role grounding in its prompt). Using multiple smaller tasks and reminding the agent of its role each time reduces this risk. If the execution agent repeatedly fails or goes off track, the orchestrator might swap in a different model (perhaps a larger model for tricky code) or escalate to human intervention.

- **Quality Control (Reviewer)** – **Responsibility:** Evaluating outputs for correctness and adherence to requirements. This agent does not produce new code (except perhaps suggestions); it acts as a *critic or tester*. In code contexts, it might perform static code analysis, spot bugs or suboptimal implementations, and ensure style guidelines are followed. It might also review documents for completeness or answers for factual accuracy, depending on the app domain. **Success metrics:** If the QC agent catches defects that would have caused failures, it's doing its job. A metric could be **defect recall rate** – e.g., of all issues later identified, what fraction did the QC agent flag? Another metric: false positives should be minimal (an overly strict QC that flags correct code can waste time). **Failure modes:** The QC agent might miss critical bugs (false negatives) or raise incorrect issues (false positives). It might also exceed its authority – e.g., trying to "fix" code during review instead of just reporting (which could introduce new errors). For reliability, some workflows explicitly separate the *detection* of issues (QC agent) from *fixing* (execution agent) to maintain clear accountability and traceability. The QC agent's findings feed back into either a rework loop or, if minor, into a task list for a later polish phase.

- **Testing / Validation** – **Responsibility:** Rigorously testing the system's functionality and verifying that all acceptance criteria are met. This role often overlaps with QC, but in our taxonomy we treat it as executing dynamic tests (running code, using the application as a user would, etc.) and performing final validations. It could be one agent or a set of agents (for different test types). For example, one agent might generate and run unit tests for each function, another might simulate user scenarios (integration tests), and yet another might double-check factual output or calculations.

In multi-agent academic research, one agent can even act as a "judge" of another's answers [33] , though in software dev the criteria are more objective. **Success metrics:** High test coverage and all tests passing. Also, ensuring that any discrepancies between the intended output and actual output are caught. The validation agent is successful when it gives a strong assurance (e.g. "All critical tests passed; performance is within expected bounds; output matches specification"). Some frameworks use a percentage confidence or a checklist that must be 100% ticked off [34] . **Failure modes:** A testing agent could itself be flawed – writing insufficient tests or mis-evaluating results. There's a risk of the tester sharing the same blind spots as the coder (especially if both are powered by similar models), which is why *diversity in agents* or model ensemble is beneficial [4] [35] . Another failure mode is not properly cleaning state between tests, leading to false pass/fail (for instance, a leftover variable making a test pass). The orchestrator should ensure the testing environment is isolated or reset. Lastly, if the validation agent reports success prematurely (a false sense of security), the team might deploy buggy software – hence the value of redundant checks or manual oversight on critical components.

- **Documentation / Scribe** – **Responsibility:** Creating and updating documents that capture the project's knowledge and progress. This includes updating design docs, writing usage documentation for new features, maintaining a changelog or progress journal, and preparing any necessary handoff materials. The Scribe agent ensures that if a new agent (or developer) joins later, they can quickly understand what's been done and why [7] . It also handles *knowledge management*: summarizing the session outcomes, decisions made (and their rationales), and lessons learned for continuous improvement [23] . In many workflows, this agent also manages **version control** – committing files to the repository and perhaps tagging versions or creating pull requests [24] . **Success metrics:** Clear, up-to-date documentation and successful syncing of artifacts to storage. One could measure documentation quality by completeness (does it cover all recent changes?), accuracy (no discrepancies with the code), and readability (formatted in a way that other agents can parse if needed). The scribe's output is often evaluated by whether subsequent agents run into confusion; well-structured continuity docs prevent redundant questions. **Failure modes:** Documentation could be stale or incorrect if the agent misinterpreted what happened. There's also a risk of verbosity – including too much irrelevant detail such that important info is lost in noise (agents operate better with concise context [22] ). Another failure is if the agent pushes unfinished or broken code to the repo, which underscores why commit actions should occur only after validation gates. Ideally, the orchestrator confirms with the validation agent that everything is good before instructing the doc agent to commit. A process with periodic backups (and perhaps a rolling undo plan) mitigates catastrophic failures (the user's guidelines explicitly mention periodic GitHub sync for recovery [24] ).

**Inter-agent Communication:** Agents communicate either **directly** (passing messages or files) or via the orchestrator as an intermediary. Common patterns include an orchestrator prompting one agent with another agent's output or multiple agents reading/writing to a shared artifact (e.g., a "planning document" on disk [36] [37] ). For example, a Claude-based system had four agents share a `MULTI_AGENT_PLAN.md` file to stay synchronized [38] [37] . Communication protocols should be explicit to avoid confusion: e.g., using structured messages like JSON for agent-to-agent handoffs or predefined file formats. Some frameworks allow agents to treat others as tools – effectively calling them via function interfaces [39] . A crucial aspect is setting **boundaries**: an agent should know what info to consume and produce, and not stray beyond that (for instance, the Builder agent shouldn't directly modify the plan document – that's the Architect's domain). Embedding role names in prompts (e.g., "You are Agent 3 – Validator. Your job is X…") helps reinforce these boundaries [40] [41] . Communication also includes *validation gates*: after an agent finishes, the orchestrator

might send its output to a "gatekeeper" (like the QC agent) with an instruction like "verify this JSON is correct and meets criteria". Only on approval does the workflow continue.

**Validation Gates & Dependencies:** Dependencies between agents create a directed flow (often a DAG – Directed Acyclic Graph – of tasks). For instance, the Execution agent should only start coding after receiving a plan JSON from the Strategist; the Tester should only run after code is generated, etc. Graph-based orchestration can capture these dependencies explicitly. Modern orchestration libraries let you define such **DAGs of agent tasks**, enabling concurrent execution where possible and enforcing order where needed. Research on DAG-based agent reasoning shows significant speedups by running independent subtasks in parallel while maintaining logical constraints [42] . In practice, this means if you have two unrelated features to implement, you might spin up two coder agents in parallel, then have a single tester combine their outputs. However, parallelism adds complexity in communication (agents might need to merge changes). A safer default is sequential with checkpoints, which can later be optimized if needed.

To summarize, every agent role has a **clear contract**: inputs it relies on, processing it does, outputs it promises. Keeping roles modular and loosely coupled (through well-defined interfaces) yields a system where components can be tested or replaced independently – akin to microservices architecture, but for cognitive tasks [9] . The **Agent Role Taxonomy** (see `taxonomy.md` ) provides a structured breakdown of each role's duties, success metrics, and failure modes for quick reference.

## Frameworks and Tools for Agent Orchestration

The rise of multi-agent systems has spurred a plethora of frameworks and libraries to manage agent orchestration. Each takes a slightly different approach to coordination, communication, and execution flow. Here we compare some leading frameworks and methodologies as of 2024–2025, including LangGraph, CrewAI, AutoGen (and its successors), Hugging Face's new "Swarms" tools, and custom DAG-based systems.

- **LangChain + LangGraph:** *LangChain* has been a popular library for building LLM applications with chains and tools; **LangGraph** is a 2024 extension focused on multi-agent workflows [43] . LangGraph represents agent interactions as a graph rather than a simple sequence, allowing complex branching, looping, and hierarchical control flows [43] . Each node in the graph is an agent or tool call, which can maintain its own state and memory. Notably, LangGraph emphasizes **fault tolerance**: nodes can retry on failure, and the graph can include error-handling paths [44] . It's designed for **stateful** operations, meaning an agent can carry context over multiple turns, and for **durability**, meaning it can recover from intermediate errors without losing the whole process. LangGraph integrates tightly with LangChain (using its abstractions for memory, tools, and prompt templates), but gives developers explicit control over the orchestration logic that was often hidden. As the LangChain team noted, LangGraph aims to make agents "more transparent and controllable – explicitly allowing developers to decide each step, each piece of context, and each tool invocation" instead of burying logic inside the model's chain-of-thought [45] [46] . In essence, LangGraph offers a robust framework for building **custom agent pipelines**: you define the DAG of which agent outputs go to which agent inputs, and the framework handles execution with concurrency and memory management. This is very useful in enterprise contexts that need **memory, fault tolerance, and adaptive retries** at scale [2] . LangGraph is open-source (part of LangChain's ecosystem) and is often used when you want fine-grained orchestration but still leverage LangChain's modules for the agents themselves.

- **CrewAI: CrewAI** is a lean, Python-first framework purpose-built for multi-agent automation. It allows you to define a "crew" of agents, each with specific roles and tools, and orchestrates their collaboration through a simple interface. CrewAI prides itself on being **fast and lightweight** – it's built from scratch (not on LangChain), which means fewer dependencies and often better performance in production scenarios [47]. The design encourages *intuitive abstractions*: you define agents and their capabilities, then specify how they should work together on a task. Under the hood, CrewAI's runtime handles the message passing and can incorporate human inputs if needed. CrewAI has been used in real-world projects; its focus on "clean code" and practical use cases resonates with engineers looking for a production-ready solution [48]. In fact, CrewAI's CEO João Moura has highlighted the framework's application in industries and even appeared alongside major AI announcements (CrewAI was featured at CES 2025 during Nvidia's keynote, underscoring its prominence [49]). **Key features:** CrewAI supports *role specialization* out of the box – you can quickly spin up agents like a Researcher, a Coder, a Critic, etc., without having to manually orchestrate each turn. It also provides hooks for *multi-turn conversations* and even hierarchical setups (manager/worker agents) [50]. CrewAI is open-source (under a permissive license) and has an active community; it's often recommended for those who want a straightforward way to deploy an agent team for automation tasks without the overhead of larger frameworks. A Reddit user summarized it as "putting together a squad of AI agents, each with their own skills, who actually get stuff done, not just talk" [51] [52] – emphasizing its action-oriented design.

- **Microsoft AutoGen and AG2:** *AutoGen* (originally from Microsoft) is a framework for building **conversational multi-agent systems** [53]. It was one of the early orchestrators that let you define multiple chat agents (backed by LLMs or even human participants) and have them interact in an automated loop. AutoGen provides high-level patterns for dialogues between agents (e.g. a "user assistant" and a "tool assistant" that cooperate), including the ability to call tools and APIs during these chats. It supports various conversation **patterns** (one-to-one, one-to-many, etc.) and emphasizes flexibility and extensibility [53]. A strong point of AutoGen is its integration with evaluation tools: Microsoft released **AutoGenBench** alongside it, which is a suite to systematically test agent performance on tasks [54]. In late 2024, Microsoft introduced **Magentic-One** built on top of AutoGen's principles, aiming at a generalist agent framework with orchestrator + specialized agents (we saw Magentic-One's orchestrator pattern earlier) [55] [56]. Meanwhile, the original AutoGen developers have spun out an open-source successor called **AG2 (Autogen2)** [49]. AG2 continues the vision of AutoGen but as an independent project, presumably incorporating lessons learned and offering more freedom (since Microsoft's version had a custom license [54]). AutoGen/AG2 are particularly strong in scenarios requiring complex **multi-turn reasoning** – e.g., agents can have lengthy discussions to plan something before acting. They also facilitate *human-in-the-loop*: e.g. an agent might defer to a human or ask for clarification via the framework if configured [53]. For developers new to multi-agent, AutoGen offered a "playground" to experiment with agents talking to each other and using tools, with an active community for support [53]. One can think of AutoGen as providing the conversational glue, while you still need to implement the actual logic of what each agent does. As frameworks like Magnetic-One and Microsoft's Agent Framework evolve, they are likely merging some of AutoGen's features (indeed, Microsoft's new Agent Framework is built on the *convergence of Semantic Kernel and AutoGen* under the hood [57]).

- **OpenAI Swarm:** In 2024, OpenAI released an educational multi-agent orchestration framework called **Swarm** [58] [59]. Swarm is a lightweight library to let developers create multiple agent instances (e.g., OpenAI GPT models with different roles) and define routines for them to cooperate.

It focuses on simplicity and approachability – providing two main concepts: *routines* (predefined interaction patterns like debate, brainstorming, or critique loops) and an event loop to manage agent turns [58] [60] . The idea is to give "a simple yet flexible way to allow agents to guide agentic collaboration" [61] . While not as full-featured as others, Swarm lowered the barrier for those who want to experiment with multi-agent strategies (it's akin to a minimal toolkit). Hugging Face's blog noted that Swarm's approach emphasizes **educational clarity** – helping users understand multi-agent concepts without a heavy framework [62] . OpenAI's Swarm also dovetails with their *Agents SDK* which provides patterns for structured tool calls and function-calling between agents. Overall, Swarm is best for smaller-scale or research use – where you might spin up a handful of agents with well-defined behaviors to see emergent outcomes. It may not have built-in persistence or error recovery like LangGraph or Microsoft's offerings, but it is highly configurable and model-agnostic. Some developers use Swarm to prototype an idea and later port to a more industrial framework once the kinks are worked out.

- **Hugging Face smolAgents:** Released at the very start of 2025, **smolAgents** is Hugging Face's entry into multi-agent frameworks [63] . The name "smol" hints at its philosophy: a *small, lightweight codebase* enabling quick creation of agents with minimal boilerplate [49] . SmolAgents' standout feature is its **code-centric approach** [64] [65] . Unlike most frameworks where agents decide actions by outputting a JSON or special string that the framework then interprets, smolAgents lets the agent directly generate **executable Python code** as its action. For example, instead of an agent producing `{"tool": "search", "args": {"query": "X"}}` , the agent would output actual Python code calling a search API function and processing the result [66] [67] . This leverages the fact that modern LLMs have been trained on a lot of code and tend to produce code reliably. According to a HF paper, this "agents writing code" approach (called CodeAct) achieved up to **20% higher success** on tool-use benchmarks compared to JSON-based formats [68] [69] . The smolAgents library provides safe wrappers: for instance, an agent can be a `CodeAgent` with certain tools (like a DuckDuckGo search function) available, and when it outputs code, that code is executed in a sandboxed environment [70] [71] . Hugging Face even integrated with an external sandbox service (E2B) to ensure that potentially dangerous code (e.g. file system writes) doesn't harm the host machine [72] . **Pros:** This approach removes an entire translation layer (no JSON->function call parsing), making execution faster and often more accurate [73] [74] . It also allows more complex logic – the agent can write a loop or conditional in code to chain multiple steps at once, which JSON schemas typically don't allow [75] . **Cons:** The flexibility means you *must trust the agent* more – it could theoretically write `os.remove("*")` and try to wipe files. SmolAgents mitigates this with sandboxing and by encouraging testing in safe environments [71] . It also requires the agent to be good at coding; with weaker code-capable models the results were not great [76] . SmolAgents is best for developers who want to prototype agent behaviors very quickly and perhaps prefer using **Python as the "glue" language** for agent logic (it appeals to those who might otherwise hand-write an orchestration script). The library is young, but given Hugging Face's ecosystem, it's poised to grow with community contributions. It's also notable that smolAgents is *LLM-agnostic* and can integrate any model from HF Hub easily [77] .

- **Custom DAG/Workflow Systems:** Beyond dedicated agent frameworks, many teams implement orchestration using general workflow tools or custom code. For instance, some use **DAG schedulers** (like Apache Airflow, Prefect, or **ZenML**) to create pipelines for LLM tasks [78] . These systems let you define steps (which could map to agent invocations or tool calls) with dependencies, and the orchestrator (scheduler) ensures the steps run in order, possibly on distributed infrastructure. The

advantage is robust scheduling, monitoring, and retrial logic out-of-the-box. ZenML, for example, integrates with LangChain and provides an interactive DAG view of runs, plus lineage tracking for artifacts [79] . However, these aren't specialized for LLM agents, so features like dynamic prompt management or long context handling must be handled by the user. Another related development is **Prompt Flow** (by Azure) – a visual tool to design prompt chains and integrate them with DevOps pipelines [80] [81] . While Prompt Flow was initially single-agent focused, it can be used to orchestrate multiple agents by calling one flow from another or by sequentially invoking prompts that represent different roles. *DAG-based parallelism* is an area of active research: as noted, approaches like **Flash-Searcher** (2025) re-imagine agent reasoning as a graph problem, enabling concurrent execution of independent reasoning paths [42] . In practice, if you have a **user-configurable DAG**, you can hand-craft an optimal workflow for your use case – for example, a search engine might have a plan where multiple query agents search in parallel on different sites, then a consolidator agent merges answers. These bespoke systems require more engineering effort but can yield highly efficient solutions tailored to a task. Companies with strict requirements (e.g. high reliability or specific infrastructure needs) often go this route: using low-level orchestration (even simple async Python scripts or state machines) to avoid the unpredictability that sometimes comes with high-level agent frameworks [82] [83] . The trend, however, is that even custom systems are adopting the best practices pioneered in the open frameworks (like schema enforcement, context management, logging hooks, etc.), essentially rolling their own "12-Factor" compliant agent orchestrators.

**When to Use What:** Choosing a framework depends on your needs. If you value **ease and quick setup**, something like CrewAI or AutoGen can get you started with minimal fuss – great for prototyping an idea or automating a personal workflow. If you require **fine control, transparency, and enterprise features** (logging, scaling, fault-tolerance), LangGraph or Microsoft's Agent Framework are more appropriate, as they emphasize giving developers control over each step and integrating with production infrastructure [84] [85] . Hugging Face's smolAgents is attractive if you want to lean into *agents as coders* and perhaps unify your agent logic and tool code in one place (Python), but you must be cautious with execution safety. Meanwhile, if your use case is novel or performance-critical, you might design a custom DAG orchestration or use research prototypes to push efficiency (for instance, a high-frequency trading app might not trust an off-the-shelf orchestrator and instead implement domain-specific agent scheduling for speed). It's also worth noting that these frameworks are not mutually exclusive – you might use LangGraph to design the high-level flow, but within one node, use smolAgents for the code generation part, etc. The ecosystem is rapidly evolving, and many frameworks are converging on similar capabilities (for example, by 2025, Microsoft's Agent Framework is bringing LangChain/LangGraph compatibility, and OpenAI's functions interface can be used within any orchestrator [86] [87] ). The key is to ensure the framework supports the **patterns you need**: be it hierarchical orchestration, parallel execution, external tool integration, or ease of debugging. In the next sections, we cover cross-cutting concerns like context management and I/O contracts that whichever framework you choose, you'll want to address.

## Context Management Strategies

Managing **context and state** across multiple agents is challenging but essential for coherent performance. Each agent typically has a limited context window (especially if using models like GPT-4 with e.g. 8K or 32K token limits). In a multi-agent setup, you cannot just dump the entire application state or conversation history into each prompt – you must be strategic about what information each agent gets and retains [88] [89] .

**Divide and Conquer:** One inherent advantage of multi-agent systems is handling *extended context* by dividing it among specialists [90] . Instead of one model trying to understand a 100-page document, you might have 5 agents each take 20 pages and then a summarizer agent combine results. By splitting context, multi-agents effectively **extend the total memory** of the system (like sharding a database). For example, in a long conversation or a large codebase, an *Information Retrieval agent* might fetch relevant pieces for the current task so that the coder agent sees only what's needed [91] . This reduces overload and keeps each prompt focused.

**Strategic Context Engineering:** It's important to **curate each agent's input** carefully. As the 12-Factor Agents principles suggest, treat context like a cache – provide only what is relevant for the current operation [88] [89] . Tactics include: - **Summarization:** Maintain a running *summary* of past interactions that is much shorter than full logs. A *rolling summary agent* can update a concise summary every N turns, which is then injected into prompts instead of raw transcripts. - **Memory Windows:** For some workflows, use a sliding window of recent content (e.g., last 2 exchanges) plus a summary of older content. This keeps immediate context and general context. - **State Feeds:** Explicitly pass state variables. Instead of expecting an agent to remember implicitly that "feature X is complete" from earlier conversation, have the orchestrator supply a state object: e.g., `state = {"featureX": "done", "outputPath": "/app/moduleX.py"}` as part of the prompt. This *structured context* reduces chances of forgetting or misremembering. - **Vector Databases:** For large knowledge bases (documentation, APIs, etc.), incorporate a retrieval step. A *knowledge agent* can do an embedding search on a vector store (like FAISS) to find the top relevant facts for the current task and provide them to the agent. This way, background info doesn't live in the main context window unless needed. - **Context Partitioning by Role:** Each role might have its own focused context. The Architect keeps the high-level spec and design decisions, the Coder has the specific function or module details, the Tester holds test cases and expected behaviors, etc. They share context only when necessary. For instance, the coder might not need the full design doc, just the part about the function it's implementing. The orchestrator can mediate: "Architect: provide the relevant design snippet for module Y to the coder agent."

**Memory Persistence:** Deciding what data to persist between sessions is important, especially if the system may shut down or agents may be restarted (which is a form of "crash recovery"). The **Continuity Document** updated by the Documentation agent is one approach – it externalizes the *cumulative memory* of the project into a file [7] . New agents can read this on startup to quickly load the state. Another approach is logging all conversations and using a *replay or warm-start* mechanism: e.g., on a new session, feed the new agent a distilled summary of the last session ("We have completed X, Y is pending, here's the plan and known issues…"). Frameworks like Microsoft's Agent Framework likely integrate with their Observability to allow *stateful agents* that can pick up context from a database or previous run logs [85] .

**Context Size and Truncation:** It's crucial to monitor token usage. If an agent's context is near the model limit, it might start dropping instructions (leading to hallucinations or omissions). Good practice is to design prompts and context such that you keep a healthy margin (e.g. aim for 50% of max tokens to be safe). Use **truncation policies**: decide what to chop off when context is too long. For example, drop the oldest conversation turns but keep the latest summary and critical data.

**Example – Planning Doc Sharing:** In the VS Code + Claude example, a *shared planning document* was used as a synchronization point [36] [38] . The Architect writes the high-level plan into `MULTI_AGENT_PLAN.md` . Then when the Builder agent starts, it reads that file (the orchestrator literally includes "Read the plan to get up to speed" in its prompt) [92] . This is a neat way to ensure the builder only sees *processed, important*

*context* (the plan) rather than the entire raw discussion that led to it. Similarly, the Validator might read the code file and any test results rather than the plan, because that's what's relevant to it.

**Tool-Integrated Memory:** Some frameworks allow attaching a *memory tool* to agents. For instance, an agent can call a "RecallMemory" tool with a key to fetch some stored info or call "SaveMemory" to store something. This is like an external KV memory that agents explicitly use. LangChain supports various memory classes (e.g. ConversationBufferMemory, etc.), and LangGraph would let you plug those into agent nodes. The advantage is modularizing memory – e.g., the coder agent might use a *short-term memory* that holds only recent function names, while the strategist has a *long-term memory store* for design rationale.

**Don't Overload Context:** A common mistake is trying to stuff everything into every agent's context "just in case." This can confuse the model or cause it to ignore important parts. It's better to keep prompts minimal and rely on targeted retrieval. As one engineer quipped, even if we had 100× larger context windows, we'd *still* need techniques like context compression and schema validation for reliable production use [83] [93] – feeding an LLM tons of raw info is rarely as effective as feeding it the right info.

In summary, treat context as a **first-class design concern**: plan how information flows and where it lives at each step. Use summaries, retrieval, and state passing liberally. By **"owning your context window"** and not leaving it to chance, you prevent a lot of potential errors [88] [89]. We will illustrate some context handling patterns in the reference architecture section (e.g., hierarchical context where the orchestrator maintains global state and agents maintain local state).

## I/O Contracts and JSON Schemas

To ensure determinism and reliability, agent interactions should be governed by **explicit I/O contracts**. This means defining the format and structure of inputs and outputs for each agent (or each type of agent action) – often using **JSON schemas** or similar structured formats. Rather than exchanging free-form text instructions and responses, the system and agents agree on a constrained format that the agent **must** follow, effectively making the model's output a function call or a data object.

**Why Structured Output:** Structured outputs *"completely transform agent development"* by eliminating ambiguity [94]. With a JSON schema, you no longer have to guess what the agent meant – you parse the JSON. It also drastically reduces the model's tendency to go off format or hallucinate irrelevant text. If the model knows it must produce, say, an object with fields `{"error": bool, "next_tasks": [...]}`, it will tend to focus its reasoning on filling those fields sensibly. As an added benefit, the schema can inadvertently guide the agent's reasoning: by deciding what fields are required, you nudge the model to think about those aspects (e.g., a `justification` field forces it to consider *why* it chose something, often leading to more coherent choices) [95] [96].

OpenAI's function calling and libraries like Pydantic make it easy to define such schemas. For example, you might define a Pydantic model or a JSON schema for a *Plan* object with fields: `steps: list[str]`, `risks: list[str]`, etc. The strategist agent's prompt can include: "You should output a JSON that matches this schema: { ...schema definition... }". The system then parses the response. If parsing fails (meaning the agent deviated), the orchestrator can detect it and either correct or ask the agent to try again. This addresses a huge pain point where one had to do fragile regex or prompt gymnastics to get well-formatted output.

**Tool Actions as JSON:** Many agent frameworks implement tool use via a JSON (or similar) action schema. The agent doesn't literally invoke the tool; it emits an *action object* which the framework interprets. The 12-Factor guide emphasizes converting natural language intents into *schema-valid commands* that can be executed deterministically [31] [32]. For example, instead of an agent saying "Searching for latest sales figures...", it would output: `{"action": "WebSearch", "query": "latest sales figures"}`. The orchestrator sees this and *executes* the WebSearch tool with the provided query, then returns the results to the agent. This pattern (akin to function calling APIs) makes agent behavior far more predictable and debuggable [97] [32], since you have a clear record of every action and can enforce allowed actions. In our context, a coding agent might output `{"action": "write_file", "file": "moduleX.py", "contents": "def foo(): ..."};` the orchestrator writes that file. If the agent tries something outside the schema (like adding an unexpected field or an unsupported action), the system can throw an error or have a guard agent evaluate it.

**JSON Schemas for Communication:** Beyond tool calls, any agent-to-agent or phase handoff can use JSON. E.g., the Strategist produces a plan in JSON, the Execution agent expects input in that JSON format (so it knows exactly where to look for each piece of info). The QC agent might output a list of issues as JSON objects with fields like `line`, `issue_type`, `description`. A structured **test report** could include test case names, passed/failed status, error messages, etc., which the orchestrator or validation agent then reads systematically.

Designing these schemas is a bit of upfront work, but it pays off by making the **contracts explicit**. It also lets you use standard validation libraries to double-check the agent's output format. Some frameworks (e.g., Google's *Agent Development Kit* for PaLM) allow you to directly specify a schema and will ensure the model's output conforms by iterative prompting behind the scenes [98].

**Example – Schema in Prompt:** A simple illustration: Suppose we want the Strategist to output a development plan. We define a schema like:

```
{
  "phase": "string",
  "tasks": [ { "id": "string", "description": "string", "depends_on":
["string"] } ]
}
```

We then prompt: *"Provide a JSON with the following format: ... (insert schema). Do not include any extraneous text."* The strategist model, if well-behaved, will return exactly that JSON. If it doesn't, we detect and correct. In practice, models like GPT-4 are quite good at this if the schema isn't too large.

**Guardrails and Type Checking:** Using schemas enables building **guardrails** – e.g., if a field is supposed to be an integer but the model gives a string, the parser will fail and you can catch it. You can then feed back an error message ("The field X must be an integer. Please correct your output.") to the agent. This loop continues until a valid output or a retry limit is reached, greatly reducing silent failures.

**Deterministic Execution:** When each step of agent reasoning is distilled into a structured action, you achieve a sort of *deterministic control flow* around the inherently stochastic core of the LLM. The idea is to constrain the randomness to the creative parts (like generating code or text) but encapsulate it in a

predictable scaffold. As one expert summarized, an ideal LLM agent system behaves like a **pure function from input state to output state** for a given step [99] [100] – meaning if you feed the same context and prompt, you expect the same structured output (this is aided by setting a low temperature for critical formatting steps). While true determinism is impossible with LLMs (they have hidden latent states), treating outputs as if they were deterministic functions helps in testing and debugging [99] .

**Standards and Reusable Schemas:** We're beginning to see community-driven schemas for common tasks. For example, JSON schema for "FAQ answer" or "GitHub issue summarization" etc., which multiple agents or tools could reuse. Adopting or adapting existing schemas can save time. Also, maintain a library of your own schemas for your project's domain – for instance, a schema for a "feature spec", one for a "test case result", etc. Keep these under version control just like API contracts.

**Example – smolAgents vs JSON:** It's worth noting the contrast: smolAgents deliberately avoids JSON by going to code, but ironically, even in that model, there's an implicit schema – the function signatures of the tools. The agent must produce code calling `generate_image(prompt)` and then `display(image)` [101] [102] ; if it doesn't follow that API, the code will error out. So either way, having a strict definition of what constitutes a valid action is crucial. JSON schemas are just a very LLM-friendly way to do this (since they align with how models were trained to output JSON often).

In our deliverables, we provide **JSON schema templates** (`schema.json`) for some typical agent communications: task list format, handoff plan format, issue list format, etc. These can serve as a starting point to tailor for specific projects.

## Safe and Deterministic Execution

Autonomous agents bring power but also unpredictability. Ensuring **safety** and **determinism** in their execution is a top priority, especially as we integrate them with real-world systems (file systems, APIs, etc.). Here we outline methods to keep agent actions controlled and reproducible:

- **Sandboxing and Environment Isolation:** Never run agent-generated code directly on a production machine without safeguards. Use a **sandbox** or isolated interpreter where the agent's code can execute with limited permissions. As mentioned, Hugging Face's smolAgents uses a sandbox (via E2B) to safely execute Python tool calls [72] . This prevents catastrophic outcomes like the agent running `rm -rf /` on your actual environment. Docker containers or virtualization can achieve similar isolation. For web browsing or external actions, consider using headless browsers or read-only modes so the agent can't, say, make unauthorized purchases or alter data. A real example: an agent instructed to "clean the repo" might try to delete files – if sandboxed, you can review those deletions before applying to the real repo. *Filesystem virtualization* (having the agent operate on a copy of the directory, then diff and apply changes after review) is a practical approach.

- **Permissioning and Policy Enforcement:** Borrowing from security, implement a **capability model** for agents. Each agent (or each action) gets a certain scope of permission. For instance, the Builder agent may only be allowed to write to a `workspace/` directory, not anywhere else. The Tester agent might be permitted to execute test scripts but not to call external APIs (unless that's needed). Define these policies clearly. Some frameworks let you set this declaratively. If an agent tries to do something outside its permission (like the coder agent attempts an HTTP request when only the

researcher is allowed internet), the orchestrator should block it. Microsoft's Agent Framework likely integrates such governance and compliance rules (given it touts compliance as a feature) [85] .

· **Deterministic Tools & Side-Effect Control:** Aim for **deterministic side-effects**. If an agent calls a tool, that tool should behave predictably. For example, if there's a tool to fetch the current time or random number, realize that introduces nondeterminism in the workflow – avoid unless necessary. If an agent reads a file listing (which could change between runs), consider sorting it or fixing its state. Essentially, treat the multi-agent workflow as you would a pure function in functional programming: minimize external state influence. Of course, some tasks require external data (search results, etc.), but try to sandbox those too (maybe use a cached search result during development, etc., to test determinism).

· **Seeding and Temperature:** For reproducibility in development/testing, use a fixed random seed for the model if the API supports it [103] . Some LLM APIs allow specifying a seed which makes outputs deterministic across runs (for a given model snapshot). In absence of that, use a low temperature (0 or ~0.2) for critical steps to reduce randomness. Save transcripts of runs so you can replay or diff if something goes wrong.

· **Validation and Approval Steps:** Insert **safety validation agents or heuristics** before executing any critical action. For instance, after a coder agent outputs code, run a static analysis or a regex check for obviously dangerous commands (like `delete`, `format C:`, etc.). If found, require human approval or at least a second agent's approval. In an example from our content: smolAgents warns that despite safeguards, an agent could still output `"delete everything" code snippet making your PC faulty` [71] . A guardrail can be: detect the string "delete" or dangerous imports (`os.system("rm`) and block execution pending review. Another example: if an agent is going to call an external API (like spend money), maybe a policy says that action must be confirmed by a human or a specialized auditor agent.

· **Retry and Rollback:** Sometimes nondeterminism or errors will occur. The orchestrator should have a **retry mechanism** – e.g., if an agent fails to produce valid JSON, prompt it again (perhaps with a hint). If code execution fails (raises exception), decide whether to let the coder agent attempt a fix or escalate. Maintain the ability to **rollback** partial changes: e.g., if an agent halfway modified a file then crashed, you might revert to a stable state before retrying. Using version control for the workspace can help – commit when things are good, and if a step goes awry, reset to last good commit and try a different approach.

· **Logging and Reproducibility:** Determinism is bolstered by thorough logging (we expand on observability later). Log every input to each agent and the output it gave. This allows you to reconstruct the sequence of events. If a certain run produces a strange result, logs let you trace which agent thought what. Moreover, logs enable offline analysis – you could simulate the entire agent run on paper by following the log, which is a sign that the process was well-structured (if logs are gibberish or missing pieces, that indicates too much was implicit). Some teams even capture a *"workflow blueprint"* – essentially the chain of actions in a run – which can serve as a template for future runs or for test cases.

· **Combine Non-deterministic & Deterministic Approaches:** There are parts of app development where creativity is needed (design, generating code) and parts that should be cut-and-dry (running

tests, deploying). Use the AI agents for the former and traditional code for the latter whenever possible. For example, an orchestrator might be a simple Python script that calls the LLM agents for the creative bits but handles the transitions and checks in regular code (ensuring determinism at the control level). Microsoft's folks phrased it as combining *"fully non-deterministic agentic orchestrations"* with *"deterministic, repeatable agents"* where needed [104]. In practice, you might have a deterministic agent or script that always runs after coding to, say, format the code (using Black or Prettier) – it's not AI, but it's part of the agent pipeline ensuring consistency.

- **Self-checking Agents:** Another promising approach is to have agents simulate or reason about the consequences of their actions before execution. For instance, a Code agent could be prompted: "What will the output of this code be? Is there any potentially harmful effect?" (like a mental dry-run). Similarly, a plan agent could double-check that a plan's steps logically lead to the goal and nothing risky is involved. This *introspection* can catch errors in advance. Anthropic's research on Constitutional AI hints at models being able to critique and refine their responses based on rules – one could apply a similar concept for safety rules in agent behavior.

- **Human Oversight and Kill-Switches:** Especially early on, keep a human in the loop. Even if mostly observing, a human operator should have the ability to abort the process if it's going off track. For example, if an agent is stuck in a loop generating nonsense, a human can step in to adjust prompts or tasks. Design the system such that pausing and intervention is possible (e.g., maybe the orchestrator waits for a "continue" confirmation at certain checkpoints if a flag is set). Over time, as confidence in determinism grows, you can reduce the frequency of human checkpoints.

**Determinism vs Adaptability:** There's a balance to strike – too many constraints can strangle the agent's usefulness (you don't want to *overly* script it, else you lose the benefit of its reasoning). The goal is *safe autonomy*: let the agent be creative and flexible where it excels, but wrap those creative bursts in structured, verifiable steps. As one commentary put it, reliable AI systems require *"wrapping the model in protective scaffolding"* [105] [106]. Achieve determinism in the overall workflow even if individual steps are stochastic. The ideal is a system that, given the same initial conditions and random seed, will always converge to a correct outcome, or at least always flag uncertainty when it arises.

In our methodology, many of the previous sections' recommendations (JSON schemas, small tasks, validation gates, etc.) feed into determinism and safety. They are all parts of a **"defensive programming"** approach for AI agents. By expecting things to go wrong and planning for it, you dramatically increase the reliability of the final application.

## Modular and Testable Build Cycles

To harness AI agents effectively, it's crucial to adopt a **modular, iterative development cycle** – very much like agile software development, but tuned for AI's strengths and weaknesses. This approach yields intermediate deliverables that can be tested and validated before moving on, and provides natural recovery points if something fails. Key practices include:

- **Small, Self-Contained Modules:** Decompose the application into modules that can be built and tested somewhat independently. This might be functional (e.g., a login component, a database module) or vertical slices. By focusing agents on one module at a time, you reduce complexity (an agent working on module A doesn't need context of module B's internals, just the interface). This

also makes it easier to pinpoint where things go wrong – if a bug appears, it's likely in the last module that was modified, not in a tangle of unrelated changes. The dev.to author explicitly noted that *splitting tasks across specialized agents yields more reliable results*, and that each piece can be optimized and tested independently [9]. For example, after finishing module A, you run all tests for A and maybe integration tests with existing parts, fix issues, and only then proceed to module B. This is analogous to human TDD or iterative dev.

- **Defined Build Phases with Gates:** Structure the cycle into phases like "Plan -> Code -> Review -> Test -> Integrate -> Document". Each phase should have a **gate condition** to pass. For instance: *Plan phase gate:* strategist and orchestrator agree on a task list with ≥95% confidence (subjective but the agent can self-assess) [12]; *Code phase gate:* code compiles and static analysis passes; *Test phase gate:* all tests green; *Document phase gate:* documentation updated and verified by diff. These gates act as **quality checkpoints**, preventing the process from moving forward on shaky ground. If a gate fails, either iterate within that phase (e.g. coder fixes issues and resubmits to tests) or, if necessary, loop back to an earlier phase (maybe the plan was flawed – go back to planning with the new info). By enforcing these gates, you effectively implement a **workflow with error recovery** at each step, rather than discovering an error only at the very end.

- **Automated Testing & Validation in-cycle:** Incorporate automated tests as part of the agent workflow, not just after. For example, after the coder agent finishes a function, an agent (or tool script) can immediately run that function's unit tests. If they fail, the system knows something's wrong before even reaching a human or production. Agents can also generate tests on the fly – a validator agent could be prompted: "Write 3 simple tests to verify the new functionality." This not only checks the code but also provides regression tests for future. Some agent frameworks allow dynamic test generation; otherwise, pre-write some test cases for expected outputs if you know them. The important point is to treat testing as *integral* to the build process, not an afterthought. One pattern is **agent delegation** for improvement: the dev.to example had a Reviewer agent that, upon finding issues, could delegate to the Writer agent to make revisions [107] [108]. This iterative refinement (like multiple edit rounds) continues until the Reviewer is satisfied. That is effectively a mini test-fix loop.

- **Continuous Integration (CI) Mindset:** Even if not literally using Jenkins or GitHub Actions, apply a CI/CD mindset. Each completed cycle (say implementing a feature) should result in a potentially shippable increment – all tests passing, code documented, integrated. Then commit that to the repo. If using a platform like GitHub, you might actually integrate with it: e.g., an agent could open a Pull Request with the changes and another agent or a human could review it. This gives an extra safety net because PRs can run CI checks. In fact, integrating agents with traditional CI pipelines is a great way to combine strengths – let the agent do the coding, then let your standard CI run lint, type-check, security scans, etc., and report back. The orchestrator agent can be configured to read CI results (e.g., via an API or email) and then decide next steps (fix or proceed).

- **Modular Prompting and Chaining:** Each agent's prompt should be modular too – don't carry an ever-growing conversation if you can break it. For instance, after a plan is made, you start a fresh context for the coder agent: it gets the plan and maybe relevant code, but not the entire conversation about how the plan was debated. This ensures that each agent's chain-of-thought is relatively short and focused (leading to better quality and less drift). Use the outputs of one phase as

inputs to the next, but you don't need to include *how* the output was obtained, just the output itself. This modular prompting helps avoid the model getting confused by irrelevant context.

- **Gate Checks and Confidence Measures:** Some gates can be model-evaluated: for example, after producing a plan, you might have the same Strategist agent do a "Planning Validation" step (which was in the user's prompt list) [109] . It reviews the plan for completeness and risks. You can implement this by literally asking the agent after it presents a plan: *"On a scale of 1-10, how confident are you that this plan will succeed without major issues? Are there any potential points of failure?"*. If confidence is low or issues are listed, iterate again on the plan. Similarly, after coding, ask the coding agent or a reviewer: *"Do you think this code meets the spec fully? Did you handle edge cases?"*. This self-reflection can be surprisingly effective at catching oversights. The **95% confidence threshold** repeated in the user's docs is a way to enforce that agents should really think and only proceed when they're quite sure [12] [110] .

- **Parallel Development & Integration Points:** In some scenarios, you might run two threads in parallel (e.g. front-end and back-end development with separate agent teams). That's fine as long as you have defined integration points where things merge and get tested together. Use integration tests at that point to catch any mismatches. If you had multiple orchestrators (one per sub-team), you might have a higher-level orchestrator or a final integration agent that pulls everything together. The key is still modular: each sub-team's output is validated in isolation, then the combination is validated. Multi-agent orchestration frameworks can handle multi-team by either running separate graphs that converge, or one big graph with branches.

- **Progressive Complexity:** Start with simple versions and iterate to add complexity. Agents sometimes hallucinate complex solutions when a simple one would do (because the prompt might make the task seem grand). A good tactic is to have the strategist outline a **minimal viable product (MVP)** first. Then potentially have an iteration where the strategist (or an enhancer agent) says "OK, now here are some enhancements to make." This iterative enhancement mirrors agile sprints. It's easier to debug a simple base and then add features than to build a giant complex thing in one go. Agents too will produce more reliable output for smaller increments.

- **Recovery Systems:** In modular cycles, **recovery** from a failure is usually limited in scope. If the tester finds a bug, you only need to go back to coding for that module. If the whole plan was wrong, you might scrap it and re-plan, but you haven't deployed anything yet – that's fine. More catastrophic is if something goes wrong in integration or deployment; to recover, you have version control to revert to a stable version. Always ensure you have a recent stable commit (the documentation agent's commit step helps here). Another angle is using *checkpointing*: after each major phase, snapshot the state (e.g., compress the working directory or save key outputs). Then if a later phase corrupts something (say an agent overwrote a config incorrectly), you can restore the snapshot and try a different approach. Recovery also involves handling agent crashes – e.g., if an LLM call fails or times out. Orchestrators should catch exceptions and either retry or spawn a fresh agent instance to continue (loading context from documentation as needed). The continuity docs are essentially there to let a new agent recover context within 15 minutes and continue the build [111] [112] . Design your process so that no single agent carries crucial state only in its "head" – always externalize to files or shared memory at checkpoints so another agent can pick it up if needed.

- **Continuous Improvement:** Each cycle provides data to refine prompts and strategies. Perhaps you notice the coder agent often forgets to update the docs string of functions; you can modify its prompt template to remind it explicitly. Or the tester agent maybe wasn't thorough, so next time you ask it for more test cases. Maintain a **living checklist** of such improvements (the lessons-learned doc). Over time, this forms a playbook that further **determinizes** the process (because you're removing sources of error).

By designing the build process to be modular and testable at every step, you end up with a **pipeline that is resilient**. It's akin to having multiple circuit breakers – a failure trips one, but doesn't bring down the whole system. This structure also builds trust: team members (or stakeholders) can see each intermediate result, verify it, and be confident moving forward. In agent-land, where things can feel like a blur of magic, such grounded checkpoints are invaluable.

## Integration with Development Tools and Platforms

To maximize productivity, autonomous agents should be integrated into the same ecosystem developers use – source control, IDEs, project management tools, etc. This not only keeps a human-friendly view of progress but also lets agents leverage existing workflows. Here are best practices for integration with **GitHub, VS Code, and orchestrator models**:

- **Version Control (GitHub) Integration:** Treat the agent team as you would a human team in terms of code management. That means using Git for tracking changes, branches for features, and pull requests for merges when appropriate. An agent (or a script supervised by the orchestrator) can automatically commit changes after each successful cycle [24] . Include meaningful commit messages (agents can be prompted to write a summary of what was done as the commit message). Pushing to a remote repository (GitHub or GitLab) after major milestones provides an off-site backup and a timeline of progress. It also opens the door to **human oversight**: a developer can review the agent's commits on GitHub and intervene if something looks off. Some have set up workflows where an agent creates a PR and a human has to approve it – a nice safety check for critical code. Additionally, by integrating with GitHub, you can use GitHub Actions to run CI tests on the agent's code. The results of those actions could then feed back into the agent system (via an API or webhook to the orchestrator, indicating success or failure). Overall, using Git imposes a helpful discipline on the agent: code changes must be consolidated, conflicts must be managed (imagine two agents changed the same file – a merge conflict might arise; it's tricky for agents to solve, so better to orchestrate so that doesn't happen, or have a strategy for it). Frequent commits (with the ability to revert) implement a **"checkpoint and rollback"** mechanism.

- **IDE Integration (VS Code, etc.):** Developers often work in IDEs like VS Code, and there's a trend to bring agent assistance directly into those environments. VS Code's *Chat Agent mode* allows you to give high-level tasks in natural language and let an AI agent autonomously carry them out in the editor [113] . There are also extensions (AI Toolkit, etc.) that offer an "Agent Builder" interface within VS Code [114] . These essentially embed an orchestrator that can manipulate your workspace. For instance, an agent can open files, edit code, run terminals – all through VS Code APIs (with your permission). Integrating your agent system with VS Code means you can watch what agents do in real-time (as if a ghost is typing in your editor). It also means you can easily step in. Several "agentic" VS Code extensions exist, like **CoderChat** or **Blackbox AI Agent**, though at the time of writing they vary in autonomy level [115] . A particularly illustrative case is a Reddit user who launched 4 Claude

agents in separate VSCode terminals, each role-bound (Architect, Builder, Validator, Scribe) [36] [37] . They coordinated via files in the workspace. This setup demonstrates that you can use a standard development environment as the canvas for multi-agent work. The agents essentially perform as specialized team members each operating in their own console/editor space. For integration, it might be as simple as scripts to spawn terminal sessions with specific prompts (like they did with `claude > You are Agent 1 - The Architect...` in each terminal) [116] . Another approach: use VS Code's Live Share or similar to have an agent join a collab session. However, note that giving agents direct write access in your IDE can be risky; many prefer to run agents on a separate branch or a forked repo and then manually merge.

- **Project Management and Tickets:** If your development uses issue tracking (Jira, GitHub Issues, Linear, etc.), you can integrate agents with that as well. For example, an orchestrator could read from a kanban board which tasks are next, spawn agents to work on them, then update the tickets when done. In fact, someone made an AutoGen agent to create Linear tickets from TODO comments in code automatically [117] – a taste of how agents can connect code and PM tools. You could similarly have an agent comment on a GitHub issue with analysis or open issues when it detects a bug. This keeps the humans in the loop via the platforms they already check.

- **Orchestrator Models (Claude, GPT-4, etc.):** The choice of model for each agent is important. Often, you want your most powerful (and expensive) model to be the **Orchestrator or Strategist**, because that agent needs the best reasoning to plan and make high-level decisions. Microsoft researchers explicitly recommend using a *strong reasoning model like GPT-4* for the Orchestrator agent in complex systems [118] . GPT-4 (or the 32K context variant GPT-4-32k) is excellent for orchestration due to its ability to handle longer prompts (like entire project summaries or multiple agent outputs) and its superior reliability in following instructions. **Claude 2** (by Anthropic) is another great choice, known for its 100K context window – useful if you want to feed a lot of background documentation or have the orchestrator handle a huge conversation. Indeed, the Reddit example used Claude for all roles and it performed well [36] . The orchestrator's tasks (planning, integrating, reflecting) are akin to what these top models excel at (complex, abstract reasoning). For execution agents that produce a lot of text (like code), GPT-4 is also common, but one might opt for a cheaper model for straightforward tasks. Some frameworks let you mix: e.g., orchestrator uses GPT-4, coder uses a fine-tuned code model (like Codex or Code Llama), tester uses GPT-3.5, etc. This can reduce cost and possibly improve performance if the specialized model is better in its domain (e.g., Codex might generate code faster/cleaner than a general model sometimes). **Claude vs GPT-4o:** The term GPT-4o in the prompt likely refers to GPT-4 (OpenAI) or GPT-4 "orchestration". Anyway, both Claude and GPT-4 have been used successfully as orchestrators. One advantage of Claude is the larger context – you might prefer it if your orchestrator needs to juggle a lot of info from many agents. GPT-4 sometimes has crisper logic but with smaller context, so you have to be more selective in what you feed it.

- **Agent Orchestrators (AutoGen, etc.):** If using frameworks like AutoGen or Microsoft's Agent Framework within your dev setup, note that some offer GUI or studio tools. AutoGen has a "no-code GUI" called AutoGen Studio [119] , which might let you visually monitor agent chats and intervene. Microsoft's Azure Agent Service integrates with their portal, meaning you could manage agents on the cloud. Hugging Face is integrating agents into their web playgrounds and Inference Endpoints. Depending on your environment (local vs cloud), choose an orchestrator that fits: for local development with quick iteration, an open-source library (LangChain, CrewAI, etc.) in a Python script

or Jupyter notebook is convenient. For deploying persistent agent systems, you might consider containerizing them or using cloud agent services.

- **Collaboration and Multi-User Scenarios:** If multiple human developers are also working, integrate agents in a way that's transparent. For instance, if an agent pushes a commit, have it mention in commit message it was AI-generated (some teams use a [AI] prefix). Also ensure that if a human is editing code simultaneously, the orchestrator is aware to avoid conflicts (one approach is to "lock" files an agent is working on). VS Code Live Share with an agent could be fun – a developer could literally watch the agent type and even chat with it to correct course ("hey AI, fix that off-by-one"). While such interactive collaboration is experimental, it's a foreseeable use-case with tools like Codespaces and ChatGPT plugins for IDEs.

- **Observability in Dev Tools:** Connect logging to interfaces devs use. For example, log agent conversations to a file that's visible in the repo, or to the console output in VS Code (perhaps in a panel). Visualize the agent DAG if possible (some frameworks provide a web view). Developers will be more comfortable if they can see what the agents are "thinking" or doing. It demystifies the process and aids trust. Integrating logging into your version control (like committing a "conversation log.md" alongside code changes) could be interesting for auditing.

In short, **meet the agents where the developers are**. By integrating with GitHub and VS Code, you ensure that normal software engineering practices (code review, testing, documentation) remain in play and that agents augment rather than bypass these processes. This integration also provides **safety valves**: a failing agent can be caught by a failing test in CI; a weird code diff can be caught in a PR; a developer can jump in if needed.

## Reference Architecture Patterns

There is no one-size-fits-all architecture for multi-agent systems, but several **common patterns** have emerged. Here we describe a few, illustrated conceptually:

1. **Single Orchestrator, Linear Pipeline:** This is the pattern we've largely described in the lifecycle – one orchestrator coordinates a set sequence: *Planner -> Coder -> Tester -> Deployer*. It's akin to an assembly line. Information flows mostly in one direction (with possible feedback loops if a part fails). This pattern is straightforward and suitable when tasks are naturally sequential. For instance, developing a single feature from start to finish. The orchestrator can be a simple loop that goes through each role in order, using the output of the previous as input to the next, and looping back as needed for fixes. **Advantages:** Simplicity, easy to follow, deterministic order. **Disadvantages:** Doesn't exploit parallelism, can be slower if some steps don't actually depend on each other.

2. **Hierarchical Orchestrators (Manager/Workers):** In this pattern, you have a top-level Orchestrator (Manager agent) which may spawn or instruct multiple *sub-agents* to work concurrently or semi-independently. For example, a Manager might decompose a project into two components and tell *Coder A* to do component 1 and *Coder B* to do component 2 in parallel. It then waits for both to finish and has a *Integration agent* combine or cross-validate their outputs. This hierarchy can have multiple levels: e.g., Manager -> Team Leads -> Workers. In fact, this mirrors human org charts. Amazon researchers have discussed hierarchical orchestrators, e.g., a supervisor of supervisors, to handle complex workflows with scaling [120] . **Use case:** Large projects where different parts can be done

independently (at least for a while). **Challenges:** Requires careful design to avoid conflicts (like two agents editing the same file) and to manage communication overhead. One agent might need to summarize its result for another to use. An example could be a Manager agent orchestrating a front-end and a back-end agent to build two halves of an application, then an integration agent ensures the API contracts match.

3. **Specialist Swarm (Decentralized Collaboration):** Here multiple agents communicate without a single boss agent. They can be on equal footing or in a network topology. For instance, consider three agents: a *Brainstormer*, a *Critic*, and a *Solver*. The Brainstormer proposes solutions, the Critic finds flaws, the Solver tries to fix them, and they loop until convergence. This pattern is more free-form and *decentralized*. OpenAI's documentation or community experiments often explore such multi-agent debates or collaborations (where each agent sees others' messages and responds). Another example: an *Ideas* agent generates design ideas, a *Feasibility* agent evaluates each, and a *Decision* agent picks the best – this is a voting/negotiation style. **Advantages:** Can leverage diversity of thought and lead to robust outcomes through discussion (reducing single-agent hallucination by consensus). **Disadvantages:** Harder to control, conversation could cycle endlessly or diverge. Usually needs a termination condition (like after N rounds or if agents agree). Also, without an orchestrator, you need some mechanism to start/stop the interaction and possibly break ties. Some frameworks allow setting up such free-for-alls with specified roles and let them at it (often as research). For production, usually a bit of orchestration or at least monitoring is added to ensure it ends.

4. **Shared Memory (Blackboard) System:** This pattern involves agents that don't talk to each other directly, but via a **shared memory or blackboard**. Think of it like a whiteboard on a wall: any agent can write or read info there. Agents monitor the blackboard for items relevant to them. For example, a *Requirement* appears on the board, a Planner agent sees it, writes a Plan on the board. A Coder agent sees the Plan, writes Code on the board (or as files), a Tester agent sees new Code and writes TestResults, etc. The orchestrator in this case might just be the blackboard logic itself (ensuring atomic updates, and maybe a rule engine to trigger agents when certain info is available). Blackboard architectures were popular in AI systems historically for multi-expert collaboration. In LLM context, one could implement a blackboard as a structured data store (like a JSON doc in memory or a database). Agents then poll or are prompted with the current state of that store. For instance, the documentation file in our earlier example served as a blackboard of sorts for status [22] . **Advantages:** Loose coupling – agents don't need to be aware of who will consume their output, they just post results. Can add new agent types easily (if they find something interesting on the board, they act). **Disadvantages:** More complex to ensure everyone sees what they need and not get confused by irrelevant data. Also potential for race conditions if two agents write conflicting info (needs a conflict resolution strategy or locking).

5. **Human-in-the-Loop Hybrid:** Architecture where human and AI agents are mixed. For example, a human product manager agent sets requirements, an AI architect makes a design, an AI coder implements, a human engineer reviews the PR, etc. In orchestrator terms, treat human inputs as part of the workflow. This might mean pausing and notifying a person at certain gates. E.g., "Plan ready – waiting for human approval before coding." Or if tests fail, maybe ping a human to decide whether to proceed or adjust requirements. Many current real deployments are of this nature – AI agents doing the grunt work, humans overseeing critical decisions (for liability or quality reasons). So the architecture includes a feedback channel to humans (could be as simple as sending an email

or Slack message with a question and having a human respond which the orchestrator then feeds into the agent system).

6. **Swarm of Homogeneous Agents:** This is more theoretical for app dev, but mentionable. That's where you have many instances of a similar agent and they collectively solve a problem (like a swarm of reasoning threads that occasionally share info). There was research on evolving orchestration where multiple candidate solutions are pursued in parallel by different agents and then the best selected [121]. For writing code, one could imagine multiple coder agents each implementing an approach to a problem; then a validator agent picks the one that passes tests. This is brute-force and costly, but could be used for extremely critical code (like N-version programming but with AI, to avoid errors by consensus).

For our application-building scenario, the **hierarchical orchestrator pattern** with a central coordinator and sequential specialist agents is the primary recommendation (since it aligns with software engineering stages). But one can incorporate elements of others, like having a pair of agents debate a design within the planning stage (small decentralized collab within one node of the hierarchy).

The **Reference Architectures document** ( `reference_architecture.md` ) sketches out a few of these patterns with diagrams and descriptions, which can guide implementation. At a high level, always ensure a clear flow of information and control. Ambiguity in who does what or when often leads to agents either stepping on each other's toes or tasks falling through the cracks.

One more specific architecture example to cement ideas:

**Example Architecture – "Claude Code Quartet"** (from the Reddit case study): - Four agent roles (Architect, Builder, Validator, Scribe) run concurrently in a shared VS Code workspace [36] [37] . - Synchronization is via files: the Architect creates `PLAN.md` . - The Builder periodically reads `PLAN.md` and implements tasks, writing code files. - The Validator watches file changes (perhaps runs tests whenever a file is saved) and updates a `REPORT.md` with issues found. - The Scribe reads code/comments and writes `DOC.md` summarizing usage or updates `CHANGELOG.md` . - A simple orchestrator (could be a script or just the user manually prompting each) ensures they operate in cycles: e.g., Architect writes plan -> Builder implements first item -> Validator tests -> Builder fixes -> loop until tests pass for that item -> Scribe documents -> Architect marks item done and moves to next. This is a nice blueprint for a *collaborative agent team* working in parallel but synchronized by documents. It achieved faster progress (parallel dev + built-in checks) and clear separation of concerns [21] .

Designing a reference architecture for your use case means identifying the roles needed, deciding how they coordinate (direct vs via orchestrator vs via shared memory), and how data flows. Once that's drawn out, you can choose frameworks that support that style or implement the coordination logic yourself.

# Observability and Logging

**"If you can't measure it, you can't improve it."** Observability is critical when deploying autonomous agents – you need to monitor what they're doing, both for debugging and for trust. Here are practices for making agent workflows observable:

- **Comprehensive Logging:** Log every significant event:
- Prompts sent to each agent (perhaps truncated if huge, but key parts).
- Agent responses, whether structured or free-form.
- Tool invocations and their results.
- Internal decisions by the orchestrator (e.g., "Agent A output failed validation, retrying with Agent B").
- Errors/exceptions, with stack traces if applicable.
- Timing information (how long each step took, tokens used, cost if calling paid API).

The logs should be structured enough to parse later. JSON Lines format (each event as a JSON object) works well. This can include fields like `timestamp, agent, event_type, content, success_flag, tokens_used` etc. For human readability, also maintain a formatted log (like a markdown or plaintext transcript with annotations). In debugging sessions, it's invaluable to see the conversation history and actions taken.

- **Real-time Monitoring:** During development, tail the logs in console or use a dashboard. If using something like LangChain, their tracing utilities can show you the chain of calls. Microsoft's Agent Framework in Azure presumably has integrated monitoring (given it talks about observability) – likely you can see agent interactions in their portal [85] . If running locally, you might hook logs to an interactive viewer (even a simple Streamlit app that updates as events come in). Real-time metrics to display:
- Current step and role ("Executing: Coder agent on task 3").
- Number of iterations so far.
- Any anomalies detected (e.g., if an agent had to retry format 3 times).

- Perhaps a visualization of the workflow graph and highlights on the active node.

- **Metrics and KPIs:** Define some key performance indicators for your agent system and track them automatically:

- *Task success rate:* percentage of tasks completed without human intervention.
- *Average retries per phase:* e.g., how often do we have to generate code twice? If this creeps up, something's off.
- *Time per cycle:* how long from planning to deployment for a typical feature (could be broken down per phase).
- *Token usage per task:* to monitor cost and efficiency.
- *Defect rate:* perhaps number of bugs found in testing per 1k lines of agent-written code, etc.

These metrics can be logged or sent to monitoring services (like Prometheus/Grafana, etc.). In an enterprise, integrating with existing APM (Application Performance Monitoring) tools is wise. For instance,

logging to Azure Application Insights or Datadog so you can set alerts (maybe alert if an agent is stuck in a loop for >5 minutes or if cost per task exceeds $X).

- **Agent Output Validation Logs:** Whenever an agent's output is checked (like JSON schema validation, tests run, etc.), log the results of that check. E.g., "SchemaValidation: PASS" or the error details. This helps diagnose if failures are due to agent or system expectations.

- **Content Logging and Redaction:** If agents handle sensitive data, you may need to redact logs for privacy (or disable logging of certain content). But for development, log as much as possible, then later you can mask things like API keys or user data in logs via simple regex or by instructing agents not to print sensitive info.

- **Observability in Multi-agent context:** Traditional software logs might not capture the nuance of "why did the AI decide that?". Consider adding *trace comments* in logs from the agent's reasoning if available. For instance, if you have the agent's chain-of-thought (some frameworks let you retrieve the reasoning token by token), logging that (or at least the final rationale it provided) can be insightful. It is similar to how one might log the decision path in a heuristic algorithm.

- **Post-mortem Analysis:** When something goes wrong, you'll rely on logs to do a *post-mortem*. For example, if the final output was incorrect or a bug got through, trace back: Did the planner miss a requirement? Did the coder misinterpret spec? Did the tester not cover that case? Having logs allows this forensic analysis, after which you can refine prompts or add tests to cover the gap. If logs show that an agent hallucinated some fact and no other agent caught it, maybe that's a sign to add a fact-checking step.

- **Logging for Continuous Learning:** Over time, log data can be used to improve the system. For instance, you can analyze logs to see common misunderstandings or failure points. If a certain type of bug appears frequently in QC logs, you might incorporate a static analysis tool or a code linter agent to catch that earlier. If an agent often exceeds token limit (logs will show prompt sizes), you know to improve context management. In a way, the log becomes a dataset. Some teams have even fine-tuned models on their own agent interaction logs to get more domain-specific behavior.

- **User-Facing Observability:** If this multi-agent system is delivering an app to a user (maybe as an assistant or so), consider exposing some of the reasoning or logs to the end-user if appropriate. For example, if an agent can't complete a request, showing a user "I tried X, it failed due to Y" can be more helpful than a generic error. However, in most app dev contexts, the "user" is the developer, so that developer should have full observability.

- **Durability and Logs:** On a related note, always flush logs frequently and consider writing to a durable store (disk, cloud log aggregator) in case the process crashes. You want the logs even up to the crash point. If the orchestrator restarts and can recover, it should possibly read the last logs or continuity notes to figure out where to resume.

Galileo's blog emphasizes evaluating agent contributions which implies you need logs to do so [122]. They mention advanced analysis like counterfactual testing (disabling an agent to see impact) [123] or *Action Advancement metrics* (percentage of interactions that advanced the goal) [124] [125]. Such analyses can be done by mining logs. In a smaller scale, you might manually review a log to judge each agent's helpfulness.

At scale, you might instrument code to mark which agent's output was used vs discarded to compute those metrics automatically.

In conclusion, logging and observability are the safety net and compass for iterative improvement. **Don't treat the agent system as a black box** – shine a light on every corner of it through logging. This will not only help fix issues but also build confidence among stakeholders (they can see a trace of what the AI did, which demystifies it). We've included recommendations in the `observability_logging.md` file for setting up a logging framework for agent workflows.

## Retrospectives and Ongoing Improvement

Building with AI agents is not a fire-and-forget process. It's vital to conduct **retrospectives** – evaluate what worked, what didn't, and continuously refine the system. Because AI behavior can drift or new edge cases appear, you want an *evaluation checklist* to regularly assess and improve the development workflow itself. This is analogous to sprint retrospectives in agile, but with focus on the agent team.

Key aspects to review after each major build cycle or deployment:

- **Goal Achievement:** Did the agent team meet the user's objectives for this cycle? If not, where did the breakdown occur? For example, if the feature doesn't meet requirements, was it a planning miss or an implementation bug? Trace through logs to find the step where things deviated from expectation.

- **Quality of Outputs:** Review the code quality, design quality, and any artifacts produced (tests, docs). Are they up to project standards? If not, perhaps incorporate a style checker or provide better examples to the agents. If the documentation agent's output was lacking, maybe it needs a better prompt or more info.

- **Error Analysis:** Catalog any errors/bugs that occurred:

- If tests failed, why did they fail? Did the coder not consider that case or did the planner omit a requirement?
- Did any agent output require multiple retries? Why – was the schema unclear, or the model struggled?

- Were there any near-misses that a human caught (e.g., in code review a human noticed something the agents didn't)? Those are learning opportunities – can we teach the agent to catch that next time (via prompt or adding an agent/tool for it)?

- **Efficiency and Bottlenecks:** How long did the cycle take, and which step took the longest? Maybe the coder agent was slow due to large code generation – could break it into smaller functions next time. Or maybe the tester spent a lot of time; perhaps tests weren't automated enough. If certain agents were idling (like waiting on others often), maybe some parallelism could be introduced.

- **Communication Gaps:** Did agents misunderstand each other at any point? For example, did the coder misinterpret the plan (meaning plan description could be improved)? Did the tester

misunderstand expected behavior (meaning requirement spec was unclear)? Ensure all these miscommunications are addressed by either clarifying formats or adding explicit notes in prompts. If an agent asked for clarification or seemed confused (it might ask questions in its output), note that and adjust.

• **Agent Appropriateness:** Evaluate if each role's agent was suited to the task. Perhaps the model used for a role was too weak or too strong (overkill). E.g., if GPT-3.5 was the coder and it struggled with complexity, consider GPT-4 next time. Conversely, if a simple job used GPT-4 and incurred cost, see if a smaller model could suffice with some fine-tuning.

• **Success Metrics:** Check the metrics you defined. If action advancement was below desired (some agents not pulling weight), consider altering or removing that agent. If bug rate is above threshold, tighten QA. The retrospective should look at metrics trends over time too – hopefully things improve as you refine the system, if not, figure out why.

• **Process Adherence:** Did the agents follow the intended process? Or did we end up bypassing some steps? For example, maybe due to urgency a human jumped in and coded directly – that's okay, but note why we couldn't rely on the agent. Next time maybe plan smaller tasks so the agent could handle it. Or if an agent kept ignoring the stop condition and went on, adjust its instructions.

• **New Risks or Requirements:** Each iteration might reveal new considerations. Maybe now that part of app is done, you realize performance optimization is needed – maybe introduce a new agent role (Performance tuner) or new tests for that. Or a security concern popped up – add a static security scan agent or tool in the pipeline (like an agent to run Bandit or npm audit, etc., and handle results). The system is extensible, so use retrospective insights to evolve the pipeline.

• **Team Feedback:** If multiple people are involved with the agent system, gather their feedback. Perhaps a human reviewer found it easier to review AI code when it was well-commented – so ensure the coder agent always writes comments. Or maybe the human PM wants the plan in a certain format – adjust prompts accordingly.

• **Checklist for Next Iteration:** Synthesize all this into a concrete plan: update prompts, add/remove tools, maybe retrain a model if you have custom data, adjust the task breakdown strategy, etc. It can be helpful to maintain a **"living checklist"** of best practices that agents (or the orchestrator) should follow. For instance:

• [ ] Strategist: Always include test plan in the design.
• [ ] Coder: If using external libraries, update requirements.txt.
• [ ] Tester: Ensure to test edge cases (empty inputs, null values, etc.).
• [ ] Documentation: Include screenshots for UI features (if applicable).

Over time this checklist can be partly automated – i.e., have an agent (or a script) verify these conditions. But initially it might be a manual list a human goes through after an iteration.

The **Retrospective Checklist** (`checklist.md`) in deliverables provides a template that covers many of these points. It's meant to be used regularly (say after each sprint or major feature). By following a consistent checklist, you avoid the temptation to skip the learning process.

Finally, don't forget to celebrate successes! If the agents shaved off 50% of development time for a feature, note that and communicate it. It helps justify the approach and also identifies what worked well so you can replicate it. For example, "Agent found and fixed 3 bugs we likely would have missed – indicates the QA step is very valuable, keep it strong."

**Continuous Refinement:** The agent system can always improve as models improve and as you fine-tune prompts or even train custom models. Keep an eye on research – new techniques like *"LLM-as-a-judge"* for better evaluation [126] or *tree-of-thoughts* for reasoning might be integrated in future. But also be vigilant about changes: if the LLM API updates and the model's style changes, you may need to adjust prompts or schema (observability helps catch such drifts).

In summary, treat the AI agents as an evolving part of the team – train them, review them, and make them better with each cycle. The retrospective is where the team (human + AI) learns collectively. As one practitioner noted, *"the only way to build impressive AI experiences is to push the model to its limits and then get it right consistently – which demands careful engineering"* [127] . That careful engineering comes from iterative improvement, guided by retrospection and data.

---

*The accompanying files and templates (* `report.md` *,* `schema.json` *,* `taxonomy.md` *,* `reference_architecture.md` *,* `observability_logging.md` *,* `checklist.md` *) provide structured resources to implement the above guidelines. By following this methodology, development teams – human and AI – can collaborate in a controlled, efficient, and innovative manner to build applications with confidence.*

---

[1] [2] [3] [44] LangGraph and the Future of Multi-Agent Orchestration in AI Infrastructure | by Atul Yadav | Medium
https://atul-yadav7717.medium.com/langgraph-and-the-future-of-multi-agent-orchestration-in-ai-infrastructure-3088ea5eaed3

[4] [5] [35] [43] [48] [53] [90] Multi-agent LLMs in 2025 [+frameworks] | SuperAnnotate
https://www.superannotate.com/blog/multi-agent-llms

[6] [7] [8] [11] [12] [13] [14] [15] [16] [17] [18] [19] [20] [22] [23] [24] [29] [30] [34] [109] [110]
user_prompts_app_bldg_draft2.txt
file://file_000000000d5c622f8ca97590fc058ec0

[9] [94] [95] [96] [107] [108] Read This Before Building AI Agents: Lessons From The Trenches - DEV Community
https://dev.to/isaachagoel/read-this-before-building-ai-agents-lessons-from-the-trenches-333i

[10] [28] [31] [32] [45] [46] [82] [83] [88] [89] [93] [97] [99] [100] [105] [106] [127] 12-Factor Agents: A Blueprint for Reliable LLM Applications
https://www.ikangai.com/12-factor-agents-a-blueprint-for-reliable-llm-applications/

[21] [36] [37] [38] [40] [41] [92] [116] How I Built a Multi-Agent Orchestration System with Claude Code Complete Guide (from a nontechnical person don't mind me) : r/ClaudeAI
https://www.reddit.com/r/ClaudeAI/comments/1l11fo2/how_i_built_a_multiagent_orchestration_system/

[25] [26] [54] [55] [56] [59] [61] [118] Microsoft's new Magentic-One system directs multiple AI agents to complete user tasks | VentureBeat
https://venturebeat.com/ai/microsofts-new-magnetic-one-system-directs-multiple-ai-agents-to-complete-user-tasks

27 What are the best models for an orchestrator and planning agent?
https://www.reddit.com/r/AI_Agents/comments/1j0f3bq/what_are_the_best_models_for_an_orchestrator_and/

33 42 Flash-Searcher: Fast and Effective Web Agents via DAG-Based Parallel Execution
https://arxiv.org/html/2509.25301v1

39 91 120 122 123 124 125 Agent Roles in Dynamic Multi-Agent Workflows: Evaluation Guide
https://galileo.ai/blog/analyze-multi-agent-workflows

47 50 How to get a multi turn conversation between agents? - CrewAI
https://community.crewai.com/t/how-to-get-a-multi-turn-conversation-between-agents/3236

49 63 68 69 70 72 SmolAgents: The Latest Innovation from HuggingFace in the Multi-Agent Arena | by MJ | Bytes & Being | Medium
https://medium.com/bytes-being/smolagents-the-latest-innovation-from-huggingface-in-the-multi-agent-arena-410de9be47ec

51 52 Forget ChatGPT. CrewAI is the Future of AI Automation and Multi ...
https://www.reddit.com/r/PromptEngineering/comments/1k8k7xz/forget_chatgpt_crewai_is_the_future_of_ai/

57 84 85 86 87 104 Introducing Microsoft Agent Framework | Microsoft Azure Blog
https://azure.microsoft.com/en-us/blog/introducing-microsoft-agent-framework/

58 Multi-agent Orchestration through OpenAI's Swarm - A Hands-on ...
https://adasci.org/multi-agent-orchestration-through-openais-swarm-a-hands-on-guide/

60 A Swarm of Agents with Llama 3.2, GPT-4o... - Hugging Face
https://huggingface.co/posts/anakin87/571953975260212

62 How OpenAI's SWARM Simplifies Multi-Agent Systems - YouTube
https://www.youtube.com/watch?v=LBih635lzps

64 65 66 67 71 73 74 75 76 77 101 102 HuggingFace smolagents: The best Multi-Agent framework so far? | by Mehul Gupta | Data Science in Your Pocket | Medium
https://medium.com/data-science-in-your-pocket/huggingface-smolagents-the-best-multi-agent-framework-so-far-313178ef3c2e?source=post_page-----410de9be47ec--------------------------------

78 79 9 Best LLM Orchestration Frameworks for Agents and RAG - ZenML
https://www.zenml.io/blog/best-llm-orchestration-frameworks

80 Azure Prompt Flow: For prototyping your AI solution
https://staslebedenko.medium.com/azure-prompt-flow-for-prototyping-your-ai-solution-420ba209ddd1

81 Azure ML Prompt Flow — from zero to MLOps | by Jake - Medium
https://medium.com/expert-thinking/azure-ml-prompt-flow-from-zero-to-mlops-5a65549e72af

98 LLM agents - Agent Development Kit - Google
https://google.github.io/adk-docs/agents/llm-agents/

103 How to make AI Agents deterministic in their responses ? : r/AI_Agents
https://www.reddit.com/r/AI_Agents/comments/1iqfn9y/how_to_make_ai_agents_deterministic_in_their/

111 112 user_prompts_app_bldg_draft.txt
file://file_000000000a48622f9b1f717a44b5af16

113 Use agent mode in VS Code
https://code.visualstudio.com/docs/copilot/chat/chat-agent-mode

114  Build agents and prompts in AI Toolkit - Visual Studio Code

https://code.visualstudio.com/docs/intelligentapps/agentbuilder

115  Top Agentic AI Tools for VS Code, According to Installs

https://visualstudiomagazine.com/articles/2025/10/07/top-agentic-ai-tools-for-vs-code-according-to-installs.aspx

117  I created an Autogen Agent which "Creates Linear Issues using ...

https://www.reddit.com/r/AutoGenAI/comments/1bp1k83/i_created_an_autogen_agent_which_creates_linear/

119  microsoft/autogen: A programming framework for agentic AI - GitHub

https://github.com/microsoft/autogen

121  Collections - Hugging Face

https://huggingface.co/collections?paper=2506.15672

126  llm-as-a-judge/Awesome-LLM-as-a-judge - GitHub

https://github.com/llm-as-a-judge/Awesome-LLM-as-a-judge