.

# Build Your Personalized Research Group: A Multiagent Framework for Continual and Interactive Science Automation

Ed Li[1,*,†]    Junyu Ren[2,*]    Xintian Pan[1]    Cat Yan[3]    Chuanhao Li[1]

Dirk Bergemann[1]    Zhuoran Yang[1]

[1]Yale University    [2]University of Chicago    [3]University of Oxford

**Abstract.** The automation of scientific discovery represents a critical milestone in Artificial Intelligence (AI) research. However, existing agentic systems for science suffer from two fundamental limitations: rigid, pre-programmed workflows that cannot adapt to intermediate findings, and inadequate context management that hinders long-horizon research. We present `freephdlabor`, an open-source multiagent framework featuring *fully dynamic workflows* determined by real-time agent reasoning and a *modular architecture* enabling seamless customization — users can modify, add, or remove agents to address domain-specific requirements. The framework provides comprehensive infrastructure including *automatic context compaction*, *workspace-based communication* to prevent information degradation, *memory persistence* across sessions, and *non-blocking human intervention* mechanisms. These features collectively transform automated research from isolated, single-run attempts into *continual research programs* that build systematically on prior explorations and incorporate human feedback. By providing both the architectural principles and practical implementation for building customizable co-scientist systems, this work aims to facilitate broader adoption of automated research across scientific domains, enabling practitioners to deploy interactive multiagent systems that autonomously conduct end-to-end research — from ideation through experimentation to publication-ready manuscripts.

**Code:** github.com/ltjed/freephdlabor
**Blog:** freephdlabor.github.io

## Introduction

The automation of scientific research through artificial intelligence (AI) is a critical step toward the realization of self-improving AI. This pursuit has spurred a recent proliferation of frameworks for automating science (Ghafarollahi et al., 2024; Ghareeb et al., 2025; Gottweis et al., 2025; Hu et al., 2025; Lu et al., 2024a; Novikov et al., 2025; Schmidgall et al., 2025a,b; Tang et al., 2025; Yamada et al., 2025; Zhou et al., 2025). However, despite demonstrations of technical feasibility, adoption of these systems into scientific practitioners' daily workflows remains limited.

A primary obstacle to the adoption of current agentic systems is their dependence on fixed workflows. These systems operate on a predefined pipeline of operations, even when individual components are powered by Language Models (LMs). This imposes a predetermined sequence of steps that cannot adapt to intermediate findings or the specific requirements of diverse scientific problems. For example, such a system could not

---

\* Equal contribution. † Correspondence: {ed.li, zhuoran.yang}@yale.edu
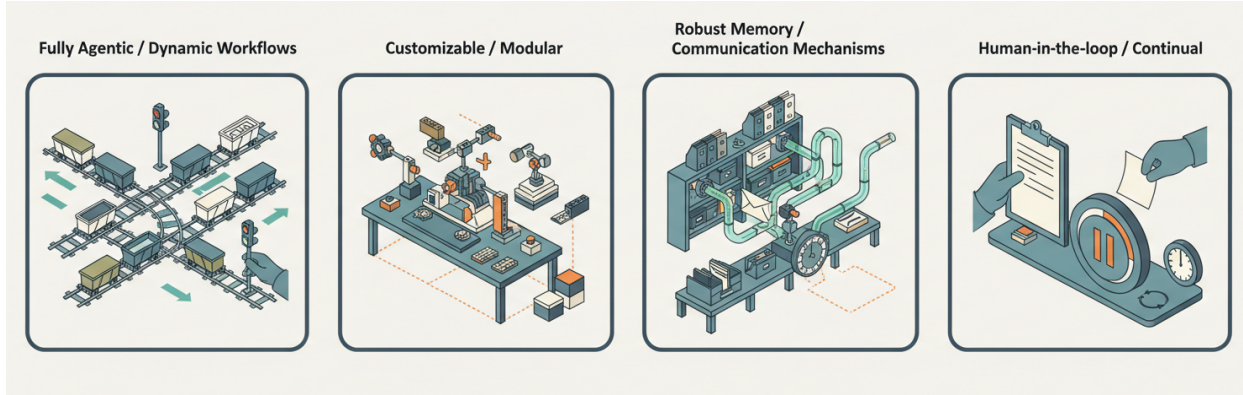
**Figure 1.** `freephdlabor` is a multiagent framework for research automation featuring: (1) dynamic workflows that adapt to real-time findings, (2) a modular architecture with customizable agents, (3) a workspace for robust communication and memory, and (4) human-in-the-loop capabilities for continual research.

pivot to a more promising research direction if initial experiments yield unexpected but valuable results. Rigid fixed workflows present a dual challenge for automated scientific research: first, they prevent the system from dynamically responding to the evolving research context, and second, their monolithic nature makes them difficult to customize for specific scientific domains without a complete architectural redesign.

In addition, automated scientific research is an inherently *long-horizon* endeavor. The process involves numerous iterative steps, including experimentation, trial-and-error, and analysis, which necessitates a large number of sequential calls to Language Models (LMs). This extended interaction horizon inevitably leads to a critical challenge in *context management*, as the volume of information can overwhelm the finite context windows of LMs. A promising approach to mitigate this is to decompose the research process using a multi-agent system. Instead of relying on a single monolithic LM, this paradigm employs multiple specialized agents, each focusing on a distinct part of the research project. For instance, a coding agent can focus solely on implementation details without being burdened by the context of manuscript preparation, thereby alleviating the context load on any single agent.

However, this multi-agent decomposition introduces new, non-trivial challenges. First, with specialized roles, agents must now communicate effectively to coordinate their actions. Second, because each agent possesses only a fraction of the total information, it must operate with only a partial observation of the project's *global state*. This gives rise to the *partial information problem*, where an individual agent may lack the comprehensive awareness needed for optimal decision-making. A third challenge lies in gracefully incorporating *human-in-the-loop* guidance; an ideal system should allow a human researcher to monitor, interrupt, and steer the research process by providing corrections or injecting domain knowledge. This paper is therefore motivated by three fundamental challenges in building effective multi-agent research systems with dynamic workflows: (i) establishing efficient inter-agent communication, (ii) maintaining a coherent view of the global research progress, and (iii) enabling seamless human-agent collaboration.

To address these challenges, we propose `freephdlabor`, an open-source multiagent framework designed for dynamic and interactive scientific discovery. Our approach is built on several core principles, each responding to the limitations discussed previously.

First, to resolve the challenge of maintaining a coherent **global state** and to counter the rigidity of fixed
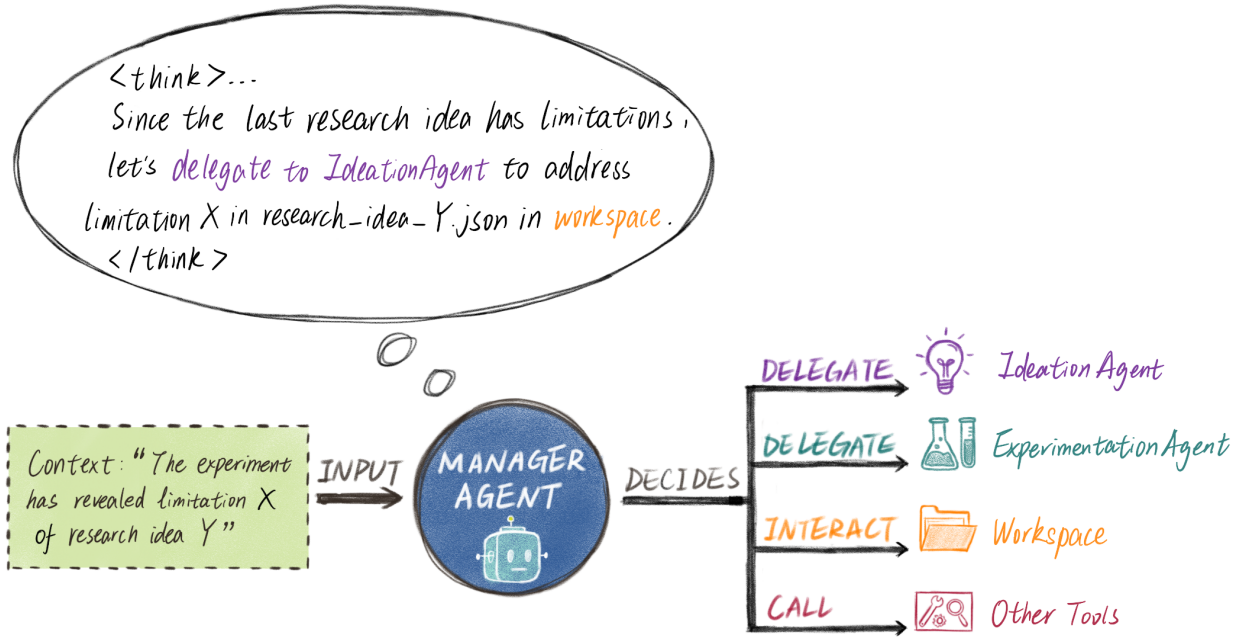
**Figure 2. Dynamic Agent Decision-Making in `freephdlabor`.** When encountering a limitation in the current research context, the system's `ManagerAgent` autonomously reasons about the appropriate response and decides whether to delegate to specialized agents, interact with the workspace, or call other tools. This dynamic decision-making enables adaptive research workflows that respond to real-time progress.

pipelines, we employ a *star-shaped architecture*. A central `ManagerAgent` acts as a coordinator, tracking overall research progress and dynamically delegating tasks to specialized agents (as shown in Figure 2). In contrast to a rigid, predefined pipeline, the `ManagerAgent` autonomously analyzes the results from the previous step to determine the most promising subsequent action, enabling the framework to adapt its strategy to a wide array of complex scientific problems. This centralized orchestration enables a *dynamic workflow* that emerges from real-time findings, rather than following a pre-programmed sequence.

Second, to establish efficient and reliable **inter-agent communication**, we directly address the information degradation inherent in long-term, language-based coordination. When agents communicate solely through string-based messages, the finite context window of LMs creates a *game of telephone* effect (illustrated in Figure 5). Over multiple conversational turns, initial instructions and critical details can be distorted or forgotten as the context becomes overloaded with intermediate reasoning. This leads to an inconsistent and lossy representation of the global state across agents (Hong et al., 2023; Li et al., 2024). To solve this, our framework implements a shared workspace that facilitates *reference-based messaging*. Instead of transcribing information, agents can directly refer to canonical data and artifacts, ensuring communication is lossless and reliable.

Third, to enable seamless **human-agent collaboration**, our framework is designed for *continual research* with integrated human oversight. The system features a real-time interruption mechanism that allows a human researcher to pause execution, provide corrective feedback, and inject domain knowledge. This, combined with memory persistence across sessions, transforms the system from a single-run tool into a collaborative partner for long-term research programs.

The primary contributions of our work are embodied in these features. The *fully agentic and dynamic workflow* provides flexibility that fixed-pipeline systems lack. The *customizable and modular architecture* allows researchers to easily modify, add, or remove agents and tools, making the system adaptable to diverse scientific domains—a concept we describe as truly plug-and-play. The framework's support for *human-in-the-loop collaboration*, via real-time interruption and guidance, transforms the system into an AI interactive partner. Finally, the *robust memory and communication mechanisms*, including automatic context compaction and the shared workspace, provide the necessary infrastructure for reliable, long-horizon research without information loss.

**Table 1. Comparison of various agentic systems for science automation.** Column "architecture" refers to whether a system entirely uses agents as the fundamental working units or partially relies on a pre-programmed chain of LM calls; column "dynamic workflow" shows whether LM outputs completely determine the flow of information in a system or not, as is the case for our system shown in Figure 2; column "customizability" refers to whether a system is modular/customizable *and* provides support features for doing so.

| | architecture | dynamic workflow | customizability | open-source |
|---|---|---|---|---|
| Agent Laboratory | fully agentic | ✗ | ✗ | ✓ |
| AI co-scientist | fully agentic | ✓ | ✗ | ✗ |
| AI Scientist | aider + LM calls | ✗ | ✗ | ✓ |
| AI Scientist-v2 | agents + LM calls | ✗ | ✗ | ✓ |
| Robin | fully agentic | ✗ | ✗ | ✓ |
| Zochi | agents + LM calls | ✗ | ✗ | ✓ |
| **freephdlabor (ours)** | fully agentic | ✓ | ✓ | ✓ |

# Related Works

Recent work on agentic systems for science has explored diverse directions, including knowledge graph-driven approaches (Ghafarollahi et al., 2024), algorithm discovery (Novikov et al., 2025), collaborative infrastructure (Schmidgall et al., 2025a), and meta-optimization (Hu et al., 2025). To provide a focused comparison, this section reviews agentic systems that target the end-to-end scientific process, from ideation and experimentation to manuscript preparation. A comparison of key features across these systems is presented in Table 1.

Early end-to-end systems often employed *hybrid architectures* that combined agentic components with structured, pre-programmed sequences of LM calls. A prominent example is Sakana AI's `AI Scientist` (Lu et al., 2024a), which gained significant attention by demonstrating the viability of end-to-end research automation. Its approach integrated a coding agent (`aider`) with a series of programmed LM calls but required a user-provided code template. The successor, `AI Scientist-v2` (Yamada et al., 2025), addressed this limitation by using a tree-search algorithm to iteratively improve code from scratch. Other systems, such as `Zochi` (Zhou et al., 2025), also adopt a partially agentic design. A common characteristic of these pioneering systems is their reliance on hybrid architectures, where agents operate within a larger, non-agentic scaffold of programmed logic.

The maturation of the agent paradigm led to the development of fully *multiagent systems*, where agentic components handle the entirety of the workflow. For example, `Agent Laboratory` (Schmidgall et al., 2025b)

orchestrates specialized agents through three fixed stages of literature review, experimentation, and report writing. Similarly, `Robin` (Ghareeb et al., 2025) achieved domain-specific success in therapeutic discovery by orchestrating three agents in a predetermined sequence. However, a noteworthy shared limitation of these systems is their reliance on human-designed, fixed workflows. The flow of information follows the same predetermined path in every run, precluding any adaptation based on intermediate findings or the evolving state of the research.

A significant step toward dynamic workflows was taken by Google's `AI co-scientist` (Gottweis et al., 2025), which runs specialized agents asynchronously based on an *a priori* allocation of computational resources. While this represents an important conceptual shift toward flexibility, its closed-source nature limits broader adoption and customization. `freephdlabor` builds on these ideas but addresses the aforementioned limitations with a distinct approach. In contrast to hybrid systems, it is fully agentic. Unlike systems with fixed pipelines, it implements a truly dynamic workflow orchestrated by a central `ManagerAgent` that makes decisions based on the real-time global state. Finally, as an open-source framework, it is designed explicitly for the customization and modularity that is necessary for broad scientific application, providing a platform for researchers to build bespoke co-scientists tailored to their specific domains.

## System Architecture

While the core features of `freephdlabor` —*dynamic workflow*, *workspace-based communication*, and *modular architecture*—are implementation-agnostic, this section details a concrete reference implementation (Figure 3). This example system, which is available for direct use or modification, serves to demonstrate the framework's design principles in action and to provide a clear blueprint for customization.

In this reference implementation, individual agents are built upon the `smolagents` library (Roucher et al., 2025) and employ the reason-then-act (`ReAct`) framework (Yao et al., 2023), which facilitates complex problem-solving through an iterative cycle of thought, action, and observation. As illustrated in Figure 4, the agent's process at each step is guided by its *memory*, a continuously updated log of its assigned task and all prior actions and their outcomes. By reasoning over this memory, the agent generates its next *action* in the form of a tool-using code segment. The execution of this code yields an *observation*—such as a tool output or an error message—which provides the agent with feedback on its action. This crucial feedback loop, where the (action, observation) pair is appended to memory, allows the agent to self-correct and refine its strategy over multiple steps until its objective is met, at which point it delivers a final response with the `final_answer()` tool. All agents in this implementation share the core `freephdlabor` features of *context compaction* and *user intervention*, which are detailed later in this report.

The behavior of each agent is fundamentally defined by its system prompt and the set of tools it can access. In `freephdlabor`, all system prompts are constructed from a unified, modular template, which provides a consistent yet flexible structure for defining agent capabilities. The structure of this template is presented below.
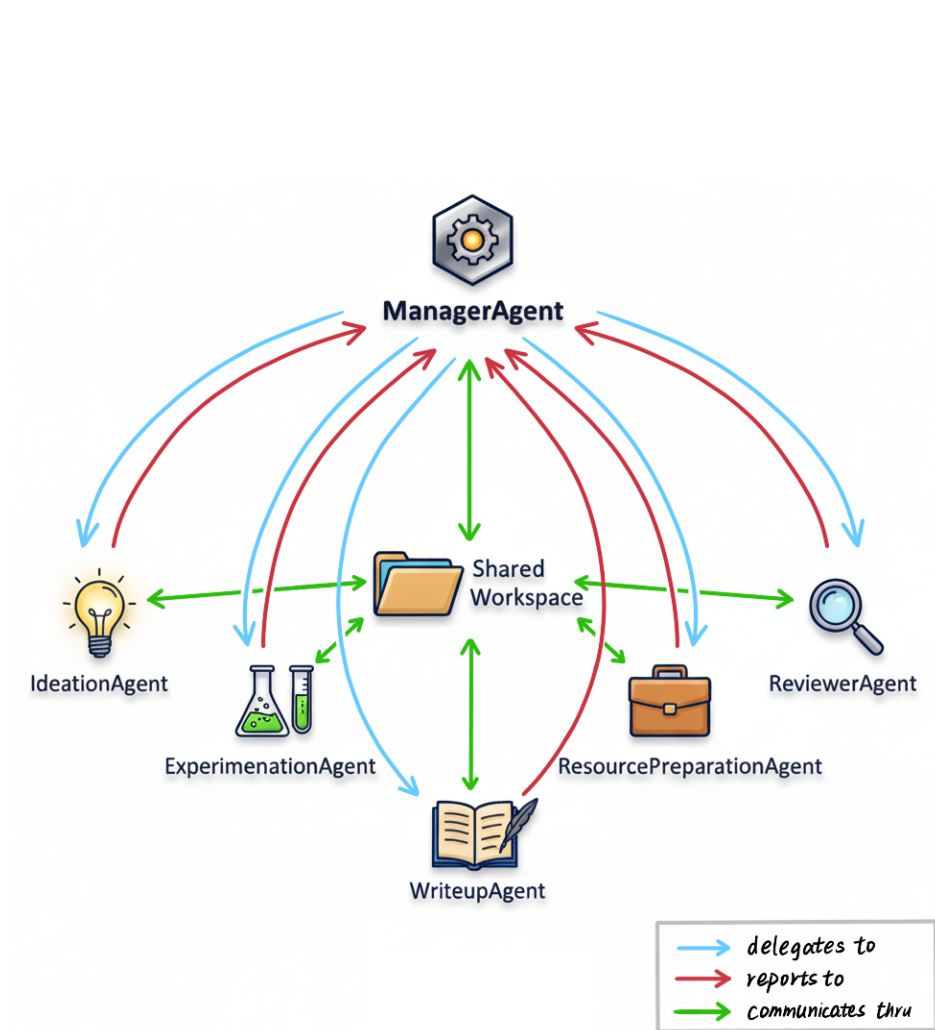
**Figure 3. Example architecture of** `freephdlabor`. Note that arrows in the figure do not indicate *workflows* like figures in other work do, but rather *options* that are available for an agent to autonomously choose from. The ManagerAgent serves as the coordinator orchestrating information flow, delegating tasks to specialized agents and managing communication through a shared workspace. All agents can read from and write to the workspace with customizable access. Thus, in addition to directly messaging each other, they may refer to files in the workspace when communicating with other agents to avoid the *game of telephone* as illustrated in Figure 5. It is important to note that this set of agents shown here is not stationary, users can modify, add, or remove agents as needed.

You are a specialized agent in a multi-agent system designed for autonomous, end-to-end AI/ML research. Your primary function is to write code blobs to call tools to accomplish given tasks. You are also given access to a workspace, which is a folder with files potentially relevant to the task at hand.

Write Python code to use the available tools and accomplish the task. During execution, you can use print() to save whatever important information you will need in memory. In the end you have to return a final answer using the final_answer() tool.

Here are a few examples using notional tools:
—
Task: "What is the result of the following operation: 5 + 3 + 1294.678?"

Thought: I will use python code to compute the result of the operation and then return the final answer using the final_answer tool
result = 5 + 3 + 1294.678
final\_answer(result)
—

... [5 additional notional tool-use examples demonstrating iterative thought→code→observation patterns including multi-step tool usage, web search with query refinement, comparative analysis, and cross-source verification]

On top of performing computations in the Python code snippets that you create, you only have access to these tools:

```
<LIST_OF_TOOLS>
[Available tools to the agent: name, description, expected input/output types for each tool.]
</LIST_OF_TOOLS>
```

Here are the rules you should always follow to solve your task:

1. Write code in ```python blocks ending with '```' to use tools and accomplish the task.
2. Use only variables that you have defined!

... [7 additional rules including constraints on tool-calling, variable naming, authorized imports, state persistence, and string syntax requirements]

ALWAYS use the correct markdown format shown in all examples above: ```python your_code_here ```

## Workspace Management

```
<WORKSPACE_GUIDELINES>
[Instructions on using the shared, file-based workspace for communication and external memory.
Also includes details about workspace tools, inter-agent communication patterns, etc.]
</WORKSPACE_GUIDELINES>
```

## Agent Instructions

```
<AGENT_INSTRUCTIONS>
[Agent-specific behavioral instructions: role definition, capabilities, workflow methodologies,
quality standards, tool usage guidance]
</AGENT_INSTRUCTIONS>


<MANAGED_AGENTS>
[Optional section only for agents with managed subagents (i.e., ManagerAgent in our example
system). Contains instructions about calling managed subagents as tools with task descriptions
and relevant context as parameters.  Lists all managed subagents with their names and
descriptions.]
</MANAGED_AGENTS>


Now Begin!
```

This compositional approach to prompt engineering allows for both structured, predictable behavior and a high degree of specialization. The template begins by instructing the agent on its basic operational pattern. In our implementation, following the convention of the `smolagents` library, agents call tools by generating Python code snippets. This design choice allows for complex logic and data manipulation to be expressed directly, though other tool-calling formats, such as JSON, could also be supported. The prompt is composed of four main modular sections, which are dynamically filled based on the agent's role:

- **`<LIST_OF_TOOLS>`:** This section is populated with the specifications for all tools available to the agent. It includes not only universally shared tools for file management but also specialized tools unique to the agent's role (e.g., the `RunExperimentTool` for the `ExperimentationAgent`). For customization, a user would typically add or define new tools here to equip an agent for a specific scientific domain.

- **`<WORKSPACE_GUIDELINES>`:** This component is identical for all agents and provides the common protocol for interacting with the shared file-based workspace. It outlines the rules for communication and collaboration, ensuring that all agents adhere to the same standards. This section is part of the core framework and is generally not modified by the user.

- **`<AGENT_INSTRUCTIONS>`:** This is the most critical section for defining an agent's unique identity. It contains a detailed description of the agent's specific role, its core responsibilities, its expected workflow, and its quality standards. When adapting the framework to a new research problem, users will spend most of their time crafting or modifying the instructions in this section to define the desired agent specialization.

- **`<MANAGED_AGENTS>`:** This optional section is used only for agents with supervisory capabilities, such as the `ManagerAgent`. It lists the sub-agents that the supervisor can delegate tasks to, effectively treating other agents as callable tools. This component is key to creating hierarchical multi-agent structures.

Full details on the content of these sections are provided in Appendices A to D.

Building on this modular prompt structure, the following subsections detail the specialized tools that complete an agent's definition. To demonstrate how these components coalesce in practice, we then present a sample execution trace illustrating the end-to-end collaborative research process.
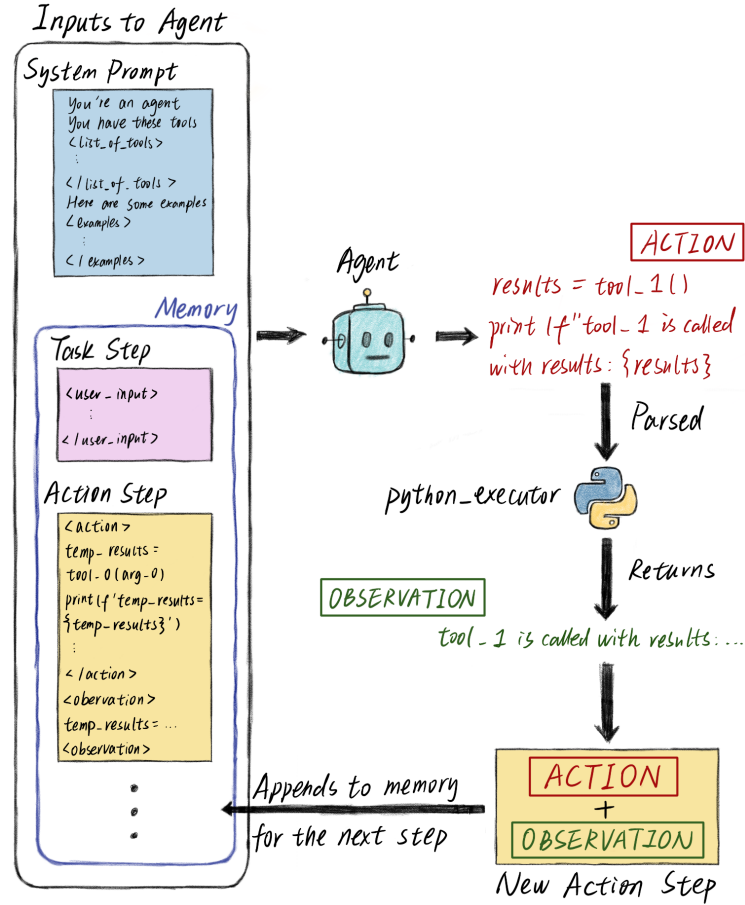
**Figure 4. Dissecting a single step of an agent.** At any given step, an agent receives inputs constructed from its own role-specific system prompt, available tools, and *memory*. *Memory* contains a task step (given by users for ManagerAgent and by ManagerAgent for other agents in our example system) and all previous action steps of this agent. It then outputs an *action*, a code blob containing tool calls, which is parsed and then executed, possibly producing some *observation* (e.g., error messages, print() statements, etc.). Then, the *action-observation* pair is appended to the agent's memory for future.

*ManagerAgent*

In a multi-agent system designed for dynamic workflows, effective coordination presents a significant scalability challenge. Ideally, each agent would need to maintain two critical types of information: (a) the entire research history (i.e., the *global state*) and (b) a description of every other agent's capabilities. However, requiring every agent to maintain such comprehensive context is unscalable, as the total information required would grow quadratically with the number of agents. To solve this, we introduce a central coordinator, the ManagerAgent, which serves as the "principal investigator"(PI) of the system. We designate the ManagerAgent as the sole agent responsible for tracking the global state and the specializations of all other agents. This *star-shaped architecture* avoids the quadratic context overhead, as each specialized agent only needs to communicate with the ManagerAgent.

Equipped with this global context, the ManagerAgent orchestrates the research workflow by invoking specialized agents as if they were functions or tools. This decision-making process is governed by the ReAct

framework (Yao et al., 2023), which consists of a two-stage cycle (illustrated in Figure 2). First, after receiving a report from a subordinate agent, the `ManagerAgent` enters a *reasoning* phase. During this phase, it does not simply pass information along; instead, it critically assesses the output, analyzing for specific success metrics, failure signals, or novel opportunities mentioned in the report. For example, it might parse a review score or identify keywords indicating a flawed experiment. Based on this analysis, the `ManagerAgent` then *acts*, selecting the most logical subsequent step. This action is what makes the workflow truly dynamic and distinct from a fixed pipeline. For instance, whereas a rigid system would be forced to proceed even with flawed results, the `ManagerAgent` can use its reasoning to break the linear sequence: it can send a paper with a poor review score back to the `WriteupAgent` for revision, or it can discard a failed experimental result and delegate a new task to the `IdeationAgent` to reformulate the core hypothesis. This capacity for contingent, context-aware routing is what allows the framework to navigate the complexities of the scientific process, making decisions that are emergent rather than pre-scripted.

### IdeationAgent

The `IdeationAgent` is a specialized agent responsible for the conceptual front-end of the research lifecycle: generating novel hypotheses through systematic literature analysis and gap identification. To do so, it is equipped with a suite of tools that enable a comprehensive methodology, mirroring the human research process of moving from broad exploration to focused synthesis and iterative refinement. The agent relies on these specialized tools to generate and refine ideas:

- **FetchArxivPapersTool**: Provides the agent with access to formal, peer-reviewed literature. By querying the arXiv API, the agent can ground its ideation process in established scientific work, identify baseline methodologies, and understand the existing state-of-the-art.

- **OpenDeepSearchTool**: Complements the formal literature search by discovering cutting-edge developments from a wide array of web sources, including blog posts, news articles, and recent, non-indexed preprints (Alzubi et al., 2025). This tool is crucial for identifying emerging trends and research gaps that may not yet be present in the peer-reviewed literature.

- **GenerateIdeaTool**: Transforms the insights gathered from the literature review into a concrete, structured research proposal. It prompts the LM to synthesize information into key fields (e.g., Name, Title, Rationale, Technical Details), ensuring that a nascent idea is articulated as a well-defined and actionable plan rather than a vague concept.

- **RefineIdeaTool**: Functions as an automated critical reviewer to iteratively improve a generated idea. The tool evaluates the proposal for logical soundness, novelty, and experimental feasibility. It can identify unjustified claims or weaknesses in the experimental design, and is also the mechanism by which feedback from downstream agents, such as the `ExperimentationAgent`, is integrated to refine the hypothesis over multiple cycles.

### ExperimentationAgent

The `ExperimentationAgent` is responsible for the empirical validation of research hypotheses. It acts as the bridge between conceptual ideas and practical results by transforming a research proposal into a functioning experiment, executing it, and processing the output. Its primary role is to manage the entire experimental workflow, from code implementation to results generation, using the following tools:

- **IdeaStandardizationTool**: Converts a research idea, which may be in a variety of natural language formats, into a standardized, machine-readable specification required by the `RunExperimentTool`. This critical pre-processing step ensures that the core technical details and experimental design of the proposed idea are accurately preserved and can be systematically implemented by the execution engine.

- **RunExperimentTool**: Executes a complete, multi-stage experimental workflow based on the standardized research idea. Adapted from the tree-search process in `AI Scientist-v2` (Yamada et al., 2025), this tool automates the process of generating and refining experimental code. Key features include: (1) **Flexible Stage Control**, which allows the agent to run partial workflows (e.g., only an initial baseline) to conserve resources, and (2) **Workspace Integration**, which saves all outputs—including code, logs, and figures—in a structured format within the shared workspace for other agents to access.

### *ResourcePreparationAgent*

A significant but often overlooked challenge in end-to-end research automation is the logistical gap between experimentation and manuscript preparation. Automated experimentation workflows, such as the one executed by the `ExperimentationAgent`, can generate hundreds of artifacts, including log files, model checkpoints, performance metrics, and deeply nested directories of plots. For a subsequent agent, like a `WriteupAgent`, navigating this complex and voluminous output is highly inefficient. It would be forced to expend a large portion of its limited context window and tool calls simply locating, parsing, and curating the correct assets, distracting from its primary task of writing.

To address this workflow bottleneck, we introduce the `ResourcePreparationAgent`, a specialized agent that acts as an intermediary data curator. Its purpose is to transform the raw, unstructured output of the experimentation phase into a clean, well-organized set of assets ready for composition into a paper. This approach exemplifies the modularity of the `freephdlabor` framework, where a dedicated agent can be inserted to handle a specific, well-defined task, thereby improving the efficiency of downstream agents.

Its core functions are enabled by the following tools:

- **ExperimentLinkerTool**: Creates a clean, accessible directory for the `WriteupAgent` by generating symbolic links to the often deeply nested experimental outputs. This abstracts away complex file hierarchies and provides a simple, flat structure for the writing process.

- **VLMDocumentAnalysisTool**: Performs deep analysis of key figures and plots from the experiment to produce high-quality textual summaries. This allows the `WriteupAgent` to understand the content of visual artifacts without needing to analyze the images directly.

- **CitationSearchTool**: Constructs a preliminary bibliography by extracting core concepts from experimental summaries and searching academic databases. It operates under strict time constraints to ensure efficiency and formats the output as clean BibTeX entries for direct inclusion in the paper.

### *WriteupAgent*

The `WriteupAgent` is an expert academic writer responsible for synthesizing all organized artifacts into a complete, publication-ready research paper. It manages the entire lifecycle of manuscript creation, from drafting individual sections to compiling the final PDF. A core design principle of this agent is its file-driven

workflow; to avoid the parsing errors common in JSON-based content exchange, all tools write LaTeX content directly to ''.tex'' files in the workspace. This ensures robustness and simplifies the generation of complex documents.

The agent's writing and compilation process is supported by a comprehensive suite of specialized LaTeX tools:

The initial draft is created section-by-section using the **LaTeXGeneratorTool**, which transforms structured experimental descriptions into formal LaTeX. To improve this draft, the **LaTeXReflectionTool** iteratively analyzes the generated ".tex" files for clarity, structure, and technical accuracy, rewriting them in place until the quality converges.

Before attempting to create a full document, the **LaTeXSyntaxCheckerTool** acts as a pre-compilation linter, identifying common LaTeX errors (e.g., unbalanced braces) and providing feedback for targeted fixes. Once the syntax is validated, the **LaTeXCompilerTool** orchestrates the final compilation. This powerful tool not only runs the LaTeX engine but also automatically resolves citations by detecting placeholders (e.g., "[cite: description]"), searching for sources with the **CitationSearchTool**, and populating the "references.bib" file.

Finally, a critical design feature is a mandatory quality gate. The WriteupAgent cannot complete its task until the **LaTeXContentVerificationTool** and **VLMDocumentAnalysisTool** validate the final PDF against a checklist of success criteria, such as adequate length, inclusion of figures, and the absence of placeholder content. This validation loop prevents common failure modes like premature termination with an incomplete or malformed paper, forcing the agent to continue its work until a high-quality output is achieved.

### *ReviewerAgent*

The ReviewerAgent acts as the system's internal quality assurance mechanism, performing the crucial function of a peer reviewer for the generated manuscript. Its primary purpose is to provide structured, critical feedback that empowers the ManagerAgent to make informed "go/no-go" decisions about the research direction. In a fully automated system, such a quality gate is essential to prevent the propagation of low-quality or erroneous results. By performing an in-depth assessment of the paper's content, methodology, and contribution, the agent generates a formal review that enables the system to either terminate a research cycle successfully or loop back for necessary revisions.

The agent relies on the following primary tool for its analysis:

- **VLMDocumentAnalysisTool**: To conduct a review that rivals a human's, the agent is equipped with a powerful vision-language model (VLM) capable of holistic document analysis. This tool performs a deep inspection of the compiled PDF, examining multiple dimensions in parallel: **linguistic quality** (grammar, clarity, coherence); **structural integrity** (logical flow, argument construction); **visual elements** (figure quality, caption accuracy); **methodological rigor** (experimental validity, statistical soundness); and **completeness** (placeholder content, missing citations). By generating a detailed, multi-faceted assessment, the tool provides the rich, nuanced information necessary for the ReviewerAgent to produce a high-quality and trustworthy peer review. This tool is also shared by other agents that require deep document understanding.

# Example Execution Trace: From Research Idea to Final Paper

To illustrate how `freephdlabor` 's architectural principles translate into practical capabilities, this section presents a summarized execution trace from a research project on "Hidden Markov Model (HMM)-based Training Phase Detection." This narrative serves as a concrete demonstration of the system's ability to handle common research eventualities, such as recovering from tool-use errors, adapting its strategy based on experimental outcomes, and iterating on a manuscript to meet a quality threshold. For clarity, we organize this continuous execution into five distinct stages. It is crucial to note that these stages were not pre-programmed; they are a post-hoc description of a workflow that *emerged* naturally from the agents' autonomous, real-time decisions.

## Stage 1: Initial Research Exploration

In a typical research workflow, initial ideation and experimentation proceed smoothly when the research direction is well-defined. Here, the system successfully generated a focused research idea and conducted initial experiments without complications. The `ManagerAgent` made standard delegation decisions, first invoking `IdeationAgent` for idea generation, then `ExperimentationAgent` for initial validation—demonstrating the expected linear progression when no obstacles arise.

---

**Stage 1**

▶ **Task Assignment:** ManagerAgent $\xrightarrow{\text{delegate}}$ IdeationAgent

*Instruction:* "Generate and refine research idea for HMM-based training phase detection"

○ **Task Execution:** IdeationAgent performed:

Web search, literature review, GenerateIdeaTool, created `working_idea.json`

◀ **Task Result:** IdeationAgent $\xrightarrow{\text{report}}$ ManagerAgent

*Response:* "Refined the idea of 'Early Warning System for Training Phases' concretely in terms of HMM methodology"

▶ **Task Assignment:** ManagerAgent $\xrightarrow{\text{delegate}}$ ExperimentationAgent

*Instruction:* "Run initial experiments on the refined research idea"

○ **Task Execution:** ExperimentationAgent performed:

Called RunExperimentTool, trained Pythia-160M on TinyStories dataset, fitted HMMs to training metrics

◀ **Task Result:** ExperimentationAgent $\xrightarrow{\text{report}}$ ManagerAgent

*Response:* "Initial experimental results completed, saved to `experiment_run/` directory"

---

## Stage 2: Workspace Configuration Error

Resource preparation typically involves straightforward directory setup and file organization. However, this stage revealed a critical integration issue: `ResourcePreparationAgent` failed to create the necessary symlink to experimental data, which `WriteupAgent` subsequently could not locate. Rather than terminating with an error, the system's dynamic workflow allowed `WriteupAgent` to attempt multiple workarounds before reporting failure back to `ManagerAgent`. This demonstrates how agents can autonomously explore solutions within their capabilities before escalating issues.

---

**Stage 2**

▶ **Task Assignment:** ManagerAgent $\xrightarrow{\text{delegate}}$ ResourcePreparationAgent ○ **Task Execution:** ResourcePreparationAgent performed:

Created `paper_workspace/` directory, analyzed the results folder to generate `structure_analysis.txt`, assembled `references.bib` for a base citable literature collection

◀ **Task Result:** ResourcePreparationAgent $\xrightarrow{\text{report}}$ ManagerAgent

*Response:* "Workspace ready but missing link to `experiment_data/`" [Warning detected]

---

▶ **Task Assignment:** ManagerAgent $\xrightarrow{\text{delegate}}$ WriteupAgent

*Instruction:* "Write research paper using organized resources from workspace"

○ **Task Execution:** WriteupAgent performed:

Searched for `experiment_data/`, attempted multiple workarounds, created placeholder `pdfs`, blocked by validation

◀ **Task Result:** WriteupAgent $\xrightarrow{\text{report}}$ ManagerAgent

*Response:* "TASK FAILED - Missing `paper_workspace/` resources, `experiment_data/` not found"

## Stage 3: Adaptive Error Recovery

In fixed-workflow systems, such failures would require human intervention to diagnose and correct. Instead, ManagerAgent autonomously analyzed the failure reports, identified the missing symlink as the root cause, and formulated a corrective strategy. By re-invoking ResourcePreparationAgent with explicit instructions to create the required symlink, then giving WriteupAgent a second attempt, the system demonstrated its ability to learn from failures and adapt its approach. The successful generation of the paper after this correction validated the effectiveness of dynamic error recovery.

> **Stage 3**
>
> ▶ **Task Assignment:** ManagerAgent $\xrightarrow{\text{delegate}}$ ResourcePreparationAgent [Corrective re-run]
>
> *Instruction:* "CRITICAL RE-RUN: Create workspace with explicit requirement for link to `experiment_data/`"
>
> ○ **Task Execution:** ResourcePreparationAgent performed:
>
> Deleted incomplete `paper_workspace/`, recreated structure, created link to `experiment_data/`, generated complete `structure_analysis.txt`
>
> ◀ **Task Result:** ResourcePreparationAgent $\xrightarrow{\text{report}}$ ManagerAgent
>
> ▶ **Task Assignment:** ManagerAgent $\xrightarrow{\text{delegate}}$ WriteupAgent [Second attempt]
>
> *Instruction:* "Write paper with corrected resources and verified data access"
>
> ○ **Task Execution:** WriteupAgent performed:
>
> Accessed `experiment_data/`, generated LaTeX sections, created figures and tables, compiled `final_paper.pdf`
>
> ◀ **Task Result:** WriteupAgent $\xrightarrow{\text{report}}$ ManagerAgent
>
> *Response:* "Paper completed successfully - `final_paper.pdf` generated"

## Stage 4: Quality Assessment and Strategic Decision

A fixed-workflow system would terminate after producing any complete paper. Here, ManagerAgent's decision to invoke ReviewerAgent demonstrates commitment to research quality over task completion. The review revealed critical weaknesses—limited scope and superficial analysis—yielding a borderline score of 5/10. At this decision point, the system could have delivered the mediocre paper to the human. Instead, ManagerAgent autonomously chose to pursue comprehensive revision, showcasing how dynamic workflows enable quality-driven decision-making without human intervention.

> **Stage 4**
>
> ▶ **Task Assignment:** ManagerAgent $\xrightarrow{\text{delegate}}$ ReviewerAgent
>
> *Instruction:* "Conduct peer review of draft paper with 5 evaluation criteria and scoring"
>
> ○ **Task Execution:** ReviewerAgent performed:
>
> VLMDocumentAnalysisTool analysis, identified weaknesses: limited scope (single model/dataset), superficial analysis
>
> ◀ **Task Result:** ReviewerAgent $\xrightarrow{\text{report}}$ ManagerAgent
>
> *Response:* "Score 5/10 - Borderline Accept, Major Revisions needed. Expand datasets, add ablation studies"

*Stage 5: Comprehensive Revision and Final Acceptance*

Following the decision to pursue quality over quick completion, this stage demonstrates the system's ability to execute complex, coordinated revisions. The workflow that emerges—expanded experimentation, workspace reorganization, comprehensive paper rewriting, and final review—showcases how agents collaborate to address systematic feedback. Each agent builds upon the previous one's work: ExperimentationAgent generates richer results, ResourcePreparationAgent organizes them effectively, WriteupAgent incorporates all improvements, and ReviewerAgent validates the quality improvement from 5/10 to 7/10. This successful revision cycle, culminating in delivery to the human, exemplifies how dynamic workflows enable autonomous systems to pursue research quality without predetermined scripts.

---

**Stage 5**

▶ **Task Assignment:** ManagerAgent $\xrightarrow{\text{delegate}}$ ExperimentationAgent [Comprehensive expansion]

*Instruction:* "Run comprehensive experiments addressing review feedback: multi-dataset and ablation studies"

○ **Task Execution:** ExperimentationAgent performed:

Expanded to multiple models (MLP, CNN), multiple datasets (IMDb, SST-2, CIFAR-10), ablation studies (features, fixed K), comprehensive diagnostics

◀ **Task Result:** ExperimentationAgent $\xrightarrow{\text{report}}$ ManagerAgent

*Response:* "Comprehensive results delivered: multi-dataset analysis, ablations, BIC analysis completed"

▶ **Task Assignment:** ManagerAgent $\xrightarrow{\text{delegate}}$ ResourcePreparationAgent [For revision]

*Instruction:* "Organize new comprehensive experimental results for paper revision"

○ **Task Execution:** ResourcePreparationAgent performed:

Created updated paper_workspace/, linked to new experiment results, generated detailed structure_analysis.txt, updated references.bib

◀ **Task Result:** ResourcePreparationAgent $\xrightarrow{\text{report}}$ ManagerAgent

*Response:* "Workspace ready for revision with comprehensive dataset links"

▶ **Task Assignment:** ManagerAgent $\xrightarrow{\text{delegate}}$ WriteupAgent [Comprehensive revision]

*Instruction:* "Write comprehensive revised paper incorporating all multi-dataset results and addressing review concerns"

○ **Task Execution:** WriteupAgent performed:

Incorporated multi-dataset results, added ablation study findings, performed deeper HMM behavior analysis, reframed narrative addressing concerns

◀ **Task Result:** WriteupAgent $\xrightarrow{\text{report}}$ ManagerAgent

*Response:* "Revised paper completed: 'On the Challenges of Detecting Training Phase Transitions...' "

▶ **Task Assignment:** ManagerAgent $\xrightarrow{\text{delegate}}$ ReviewerAgent [Final review]

*Instruction:* "Conduct final peer review of the revised paper"

○ **Task Execution:** ReviewerAgent performed:

Verified review concerns addressed, evaluated expanded scope, assessed ablation quality, confirmed major issues resolved

◀ **Task Result:** ReviewerAgent $\xrightarrow{\text{report}}$ ManagerAgent

*Response:* "Score 7/10 - Accept with Minor Revisions. Excellent job addressing concerns"

◀ **Final Deliverable:** ManagerAgent $\xrightarrow{\text{report}}$ Human

*Completion:* "Research project finished successfully. Final paper available at paper_workspace/final_paper.pdf"

---

This execution demonstrates several key aspects of `freephdlabor`'s design: (1) *Dynamic workflow adaptation*—the ManagerAgent makes real-time decisions based on agent outputs rather than following predetermined sequences; (2) *Robust error recovery*—when the missing symlink issue arises, the system identifies and corrects the problem autonomously; (3) *Quality-driven iteration*—review scores drive substantive improvements rather than premature termination; (4) *Workspace-based coordination*—agents communicate through structured files, avoiding information loss from string-based messaging; and (5) *Flexible agent invocation*—each agent makes a variable number of tool calls based on task requirements, not fixed procedures.

# Infrastructure Features of `freephdlabor`

While the architectural components described in the previous section define the fundamental organization of `freephdlabor`, their effectiveness relies critically on a set of supporting infrastructure features. These features address the practical challenges that arise in long-horizon, multi-agent research automation: managing finite context windows, ensuring reliable inter-agent communication, enabling human oversight, and preserving research progress across sessions. This section describes five core infrastructure components that collectively enable `freephdlabor` to operate reliably over extended research programs: the workspace system for robust communication, workspace tools for file-based coordination, prompt optimization mechanisms, context compaction for managing memory constraints, memory persistence for cross-session continuity, and real-time user intervention capabilities.

## Workspace System

A fundamental challenge in multi-agent coordination is the degradation of information through repeated inter-agent communication. When agents communicate solely through string-based message passing, each information exchange requires explicit transcription of data from one agent's context into a message, which is then incorporated into another agent's context. This process, repeated across multiple conversational turns, creates what we term the *game of telephone effect*: a systematic information loss analogous to the children's game where messages degrade through successive retelling.

The mechanism of this degradation operates as follows. As agents engage in extended interactions, their finite context windows become saturated with intermediate reasoning, partial results, and coordination overhead. When Agent A must communicate a complex data structure or experimental result to Agent B, it cannot pass the data directly; instead, it must serialize the information into natural language within the message. Agent B then reconstructs this information from the linguistic description, introducing potential for misinterpretation or loss of precision. Over multiple such exchanges, critical details—such as specific hyperparameter values, exact experimental configurations, or subtle patterns in results—can be distorted, omitted, or misremembered (Hong et al., 2023; Li et al., 2024). This information degradation poses a severe threat to research reliability, as downstream decisions may be based on incomplete or corrupted representations of the actual experimental state.

To address this information degradation, `freephdlabor` implements a workspace-based communication paradigm. Rather than serializing data into string messages, agents persist important information as files within a shared workspace directory. Inter-agent messages then contain only file references (paths and optional summaries), while the actual data remains in its original, canonical form. This reference-based messaging eliminates transcription errors and preserves full data fidelity. Additionally, workspace files function as persistent external memory, allowing agents to revisit prior results without relying on context-window-limited conversation history.

To maintain navigability as workspace contents accumulate over extended research runs, the framework enforces organizational structure. In our reference implementation, each agent is assigned a dedicated subdirectory, with expected structure conventions specified in the agent's system prompt to ensure consistent file placement and naming.
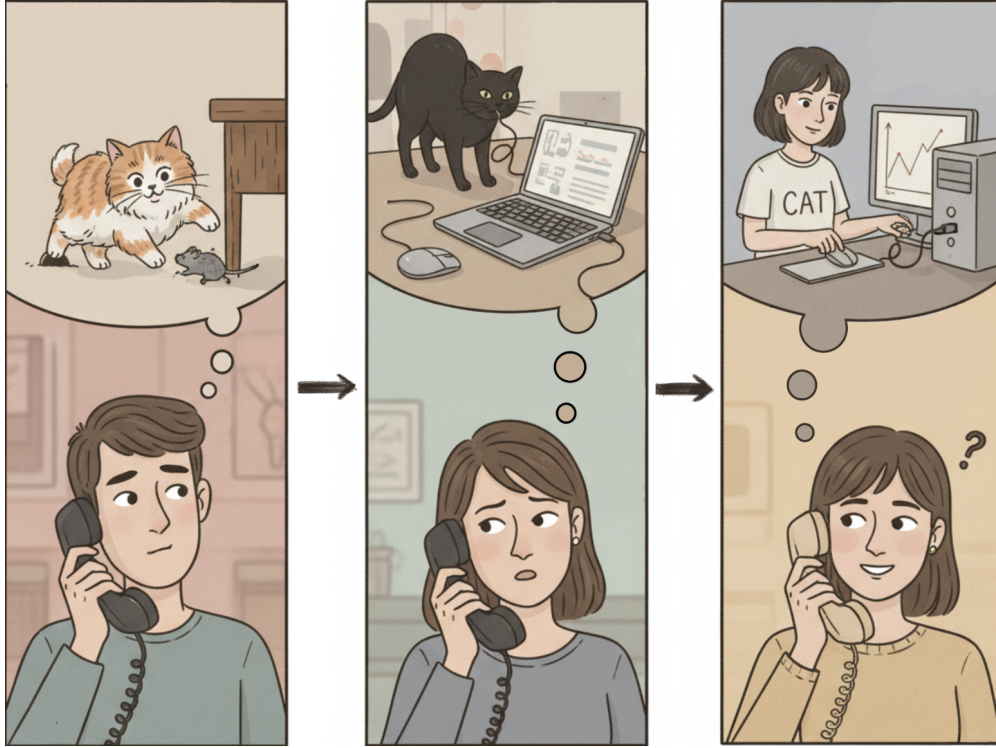
**Figure 5. Information Degradation in String-Based Inter-Agent Communication.** Visual illustration of the *game of telephone effect* in multi-agent systems. The figure depicts how information fidelity deteriorates across successive communication hops. In this example, an initial concept (a cat) undergoes progressive distortion through agent-to-agent message passing: the first agent observes the original subject, the second agent receives and reinterprets a description (resulting in a different representation), and the third agent receives further degraded information (reducing the concept to the text "CAT" on clothing). Each transcription step introduces potential for information loss, misinterpretation, or abstraction. In contrast, workspace-based communication via direct file references preserves the original data, eliminating this degradation pathway.

## Workspace Tools

To enable effective file-based communication between agents, `freephdlabor` uses a comprehensive file editing toolkit that all agents share. This toolkit consists of six core tools that abstract file operations within the secure workspace environment:

○ **SeeFile**: Reads workspace files in minimal format string for easy understanding and editing, reducing "game of telephone effect" in repeated reading and writing, optimized for code files, configurations, and text documents

○ **CreateFileWithContent**: Creates new plain text files (e.g., .txt, .py, .md) with specified content, essential for generating experiment scripts and documentation

○ **ModifyFile**: Modifies existing files by replacing specific lines, enabling precise code edits and content updates with proper indentation preservation

○ **ListDir**: Explores directory structure to understand workspace organization and locate relevant files for inter-agent coordination

- **SearchKeyword**: Searches for keywords in files or recursively within directories, returning matches with context lines for efficient information retrieval

- **DeleteFileOrFolder**: Removes files or directories when cleanup is needed, with safety constraints to prevent accidental workspace deletion

All file operations are restricted to the designated workspace directory through path validation that prevents directory traversal attacks while providing clear error messages to guide agent behavior. This security model ensures agents can collaborate freely within their allocated workspace without compromising system integrity.

## Optimizing System Prompts for Better Communications

For agents in a multiagent system to collaborate towards the same goal, we need to ensure that each agent receives the information it needs to succeed at its given task *and* faithfully communicates its work to other agents. The star-shaped architecture shown in Figure 3 already simplifies this by only requiring individual agents to communicate effectively with ManagerAgent (as opposed to all other agents).

To further simplify communication without modifying underlying LM weights (more on that in the Discussion section), we focus on improving system prompts and tool descriptions. `freephdlabor` automatically tracks all LM calls made by all agents in temporal order, creating a comprehensive interaction log. Recent literature (Agrawal et al., 2025; Zhang et al., 2025a) demonstrates that systematically examining these LM call traces, especially when collected across different runs, enables capable coding assistants or fine-tuned LMs to identify key improvement opportunities. Small fine-tuned models such as `AgenTracer-8B` show particular promise for scaling this approach (Zhang et al., 2025a). To facilitate this analysis, we provide Claude Code slash commands such as "/analyze_agent_context" and "/refine_agent_prompt" (see code for details).

## Context Compaction

A crucial goal of `freephdlabor` is to enable sustained, long-term exploration of research directions as continual *research programs*, rather than merely one-off *attempts*. This requires the system to handle two fundamental challenges: managing growing conversation context as agents reason through complex multi-step workflows, and preserving research progress across execution sessions so work can be resumed and extended over time.

Our implementation uses a callback-based automatic compaction system integrated into BaseResearchAgent. The ContextMonitoringCallback monitors memory after each ActionStep, estimating token usage through character-based heuristics (total characters divided by 4, plus tool schema overhead). When estimated tokens exceed a safety threshold—by default 75% of the model's maximum context limit—automatic compaction triggers. Compaction proceeds in three phases: (1) *External backup*: All ActionSteps to be compacted are serialized to jsonl files in workspace_dir/memory_backup/, preserving complete conversation history including tool calls, observations, model reasoning, errors, and timing information. (2) *Intelligent summarization*: The compactor extracts comprehensive context across multiple dimensions (tool usage statistics with recent call details, key observations prioritized by recency and size, recent model reasoning, encountered errors, final outputs) and generates a structured summary that preserves task continuity.

(3) *Memory reconstruction*: The agent's memory is rebuilt with one compacted `ActionStep` containing the summary plus the last 3 meaningful `ActionSteps` (those with tool calls, observations, or outputs), maintaining short-term context while dramatically reducing token count.

This approach allows theoretically unbounded conversation length while staying within model context limits. The external backup system ensures no information is permanently lost—full conversation history can be reconstructed from `jsonl` files if needed for debugging or analysis. Compaction frequency is constrained by a minimum interval (default: 3 steps between compactions) to prevent excessive summarization overhead.

## Memory Persistence and Resume

Context compaction addresses the first challenge of managing growing conversation context within a single execution session. The second challenge—preserving research progress across sessions—is addressed by our memory persistence and resume capability. Together, these two features enable the system to explore research directions as continual *research programs* rather than one-off *attempts*.

The system automatically saves the complete memory of all agents, including every execution step with detailed reasoning traces, tool usage history, and inter-agent interactions. This persistent memory captures not only high-level research progress but also the granular decision-making process that led to current findings. When combined with workspace files that serve as external memory, this creates a comprehensive record of the entire research trajectory.

When resuming a research session, the system reconstructs the entire multi-agent environment from the saved state, allowing agents to continue exactly where they left off. The user just needs to specify the workspace they wish to continue from with memory files in place. The resume mechanism enables running `freephdlabor` to explore a dedicated direction of your choice without loss of previous context.

## Real-Time Human Intervention

To enable seamless human-agent collaboration, `freephdlabor` incorporates a non-blocking interruption mechanism that balances agent autonomy with human oversight. The system continuously monitors for user intervention signals in a background process while agents execute their workflows. Unlike synchronous interruption approaches that pause execution at every step to poll for input, our asynchronous design allows agents to operate without interruption overhead until a human actively signals intent to intervene.

The mechanism is implemented via callback functions integrated into the agent execution loop. Following each action step, the callback checks for pending intervention signals. When a signal is detected, the agent suspends its current workflow and prompts the human operator for guidance, which may take the form of task refinement, corrective feedback, or initiation of a new research direction. This guidance is then incorporated into the agent's memory as a high-priority task instruction, and execution resumes with the updated objective.

This design preserves the benefits of autonomous operation while enabling precise human steering at critical junctures, transforming the system from a fully autonomous tool into an interactive research collaborator. The non-blocking architecture ensures that human intervention remains optional rather than mandatory, allowing researchers to supervise high-stakes decisions without micromanaging routine operations.

# Discussion

**Agent Deception**: Agents in `freephdlabor` can exhibit deceptive behavior under stringent requirements. For example, when the `ExperimentationAgent` is asked to produce a pdf with a length requirement, it may generate a *placeholder* document with low-information content. This mirrors broader findings on multi-agent system failures—especially task verification and inter-agent misalignment (Cemri et al., 2025)—and aligns with evidence that deceptive strategies can persist despite safety training (Hubinger et al., 2024). Multi-agent settings also introduce risks of covert coordination via steganographic channels (Kovač et al., 2024). Emerging evaluation frameworks for deception/trust (Park et al., 2025) and work on long-horizon supervisor–performer interactions (Pan et al., 2025) motivate integrating *deception checks* into our existing quality gate (e.g., `LaTeXContentVerificationTool`) and exploring a dedicated deception-auditor agent.

**Emergent vs. Pre-designed Workflows**: A line of research optimizes workflows *before* deployment via meta-search and search-over-code such as `ADAS` (Hu et al., 2025), `Darwin Gödel Machine` (Zhang et al., 2025b), `IGE` (Lu et al., 2024b), and `AFlow` (Zhang et al., 2024). In contrast, `freephdlabor` emphasizes *runtime* routing: the `ManagerAgent` reallocates work among agents/tools from real-time signals (progress, errors, user input). Surveys of LM multi-agents discuss orchestration modes and support the need for dynamic coordination beyond fixed pipelines (Guo et al., 2024a; Wang et al., 2025a). Contemporary systems illustrate differing philosophies: `PiFlow` imposes principle-aware, information-theoretic guidance to avoid aimless hypothesizing (Pu et al., 2025), whereas `freephdlabor` fosters emergent coordination without pre-programmed task sequences. Lab-in-the-loop systems like `AutoLabs` report similar benefits for reliability and self-correction (Wang et al., 2025b). Google's `AI co-scientist` demonstrates asynchronous orchestration that our runtime routing approach echoes (Gottweis et al., 2025).

**Adapting the system to individual use cases**: A natural extension of `freephdlabor` involves adapting existing agents to domain-specific requirements through tool substitution and prompt modification. For instance, for a materials scientist, substituting the `RunExperimentTool` of `ExperimentationAgent` (which is designed to run AI/ML experiments) for a tool that takes in a hypothesis and outputs experiment results. Our `RunExperimentTool` can be swapped for domain-specific executors with the same standardized I/O (idea specification $\rightarrow$ results bundle), enabling plug-and-play customization. Recent work such as `Robin` (Ghareeb et al., 2025) demonstrates that multiagent systems can achieve significant success in experimental domains, successfully identifying treatments for dry age-related macular degeneration, though such systems require substantial upfront effort to engineer fixed workflows tailored to the specific problem. `freephdlabor`'s architecture reduces this barrier by enabling adaptation through tool substitution and agent modification rather than complete workflow redesign. Resources such as `ToolUniverse` (Gao et al., 2025) provide curated collections of validated tools that can be seamlessly integrated into agent definitions. The broader landscape of tool learning with large language models (Qu et al., 2024) provides theoretical foundations for why and how to integrate external tools effectively, supporting our tool-centric *modularity* approach. Stable benchmarking frameworks (Guo et al., 2024b) ensure reliable tool use in long-running agent workflows, addressing concerns about tool reliability over extended research sessions. Domain-specific applications, such as tool-augmented agents in remote sensing platforms (Singh et al., 2024), demonstrate concrete examples of how specialized tooling translates to improved performance in targeted scientific domains.

**Improving the system via in-context learning methods**: The primary mechanism for inter-run learning

in `freephdlabor` is in-context learning. This is achieved by incorporating historical information into system prompts or by initializing the workspace with relevant artifacts from prior sessions. The structured LM calls log kept by `freephdlabor` provides trajectories for reflective prompt evolution approaches. Recent advances in reflective prompt evolution (Agrawal et al., 2025) demonstrate that systematic prompt optimization can outperform reinforcement learning approaches, directly supporting our auto-prompt optimization capabilities. Additionally, research on self-discovery mechanisms (Zhou et al., 2024) shows that models can autonomously identify missing skills through meta-prompting strategies, aligning with our vision of agents learning to improve their own capabilities between runs through structured reflection on past performance.

**Improving the system via multiagent RL**: Previous in-context learning methods also have their drawbacks. The information inserted takes up the precious context window and even distracts agents when the task is unrelated to saved information. An underexplored advantage of the multi-agent paradigm is agent-specific specialization via fine-tuning. A critical challenge in this approach is balancing the capacity constraints of post-training: augmenting individual agent capabilities without catastrophic interference to other competencies. As mentioned earlier, `freephdlabor` tracks the LM calls (i.e., state-action pairs) of different agents, serving as offline data and reward signals for collaborative RL approaches. Thus, it would be interesting to fine-tune agents using a curated version of those trajectories. Recent work on multi-agent post-co-training (Zhang et al., 2025c) demonstrates how collaborative reinforcement learning can improve multi-agent coordination by training multiple LMs with RL on collaboration signals—directly relevant to post-training our specialist agents for better cooperation. Group-relative policy optimization variants tailored for multi-agent systems (Xiang et al., 2025) provide specific algorithmic frameworks for this approach. Sequential cooperative fine-tuning methods (Lu et al., 2024c) offer templates for how agent pairs can co-evolve their capabilities over time, potentially enabling our agents to develop complementary specializations that improve overall research effectiveness.

## Conclusion

In this work, we introduced `freephdlabor`, a multi-agent framework designed to overcome the rigidity of previous systems for automated scientific discovery. Compared to prior approaches that rely on fixed, pre-programmed workflows, `freephdlabor` delegates control to the agents themselves, enabling a dynamic workflow that adapts to the real-time progress of research. This agent-centric design enables developers to focus on constructing specialized, high-quality agents while delegating coordination and integration responsibilities to the framework's support infrastructure. By allowing individual agents to be easily added, modified, or replaced, our framework provides a flexible and customizable platform for researchers across different domains. The core contribution of this work is not just a single system but a new framework for building bespoke co-scientists.

A long-term objective of this research direction is to facilitate broader access to agentic AI for science, enabling individual researchers across diverse domains to leverage automated research assistance. To this end, we provide an open-source framework that prioritizes adaptability through robust infrastructure for memory management, inter-agent communication, and human oversight. This framework offers scientists both conceptual guidelines and practical tools to construct customizable co-scientist systems tailored to their specific research needs.

# References

Agrawal, Lakshya A, Shangyin Tan, Dilara Soylu, Noah Ziems, Rishi Khare, Krista Opsahl-Ong, Arnav Singhvi, Herumb Shandilya, Michael J Ryan, Meng Jiang, Christopher Potts, Koushik Sen, Alexandros G. Dimakis, Ion Stoica, Dan Klein, Matei Zaharia, and Omar Khattab (2025). "GEPA: Reflective Prompt Evolution Can Outperform Reinforcement Learning". arXiv: 2507.19457 [cs.CL].

Alzubi, Salaheddin, Creston Brooks, Purva Chiniya, Edoardo Contente, Chiara von Gerlach, Lucas Irwin, Yihan Jiang, Arda Kaz, Windsor Nguyen, Sewoong Oh, Himanshu Tyagi, and Pramod Viswanath (2025). "Open Deep Search: Democratizing Search with Open-source Reasoning Agents". arXiv: 2503.20201 [cs.LG].

Cemri, Mert, Melissa Z. Pan, Shuyi Yang, Lakshya A. Agrawal, Bhavya Chopra, Rishabh Tiwari, Kurt Keutzer, Aditya Parameswaran, Dan Klein, Kannan Ramchandran, Matei Zaharia, Joseph E. Gonzalez, and Ion Stoica (2025). "Why Do Multi-Agent LLM Systems Fail? (MAST)". arXiv: 2503.13657 [cs.AI].

Gao, Shanghua, Richard Zhu, Pengwei Sui, Zhenglun Kong, Sufian Aldogom, Yepeng Huang, Ayush Noori, Reza Shamji, Krishna Parvataneni, Theodoros Tsiligkaridis, and Marinka Zitnik (2025). "Democratizing AI scientists using ToolUniverse". In: arXiv: 2509.23426 [cs.AI].

Ghafarollahi, Alireza and Markus J. Buehler (Sept. 2024). "SciAgents: Automating Scientific Discovery through Multi-Agent Intelligent Graph Reasoning". arXiv: 2409.05556 [cs.AI].

Ghareeb, Ali Essam, Benjamin Chang, Ludovico Mitchener, Angela Yiu, Caralyn J. Szostkiewicz, Jon M. Laurent, Muhammed T. Razzak, Andrew D. White, Michaela M. Hinks, and Samuel G. Rodriques (May 2025). "Robin: A Multi-Agent System for Automating Scientific Discovery". arXiv: 2505.13400 [cs.AI].

Gottweis, Juraj, Wei-Hung Weng, Alexander Daryin, Tao Tu, Anil Palepu, Petar Sirkovic, Artiom Myaskovsky, Felix Weissenberger, Keran Rong, Ryutaro Tanno, Khaled Saab, Dan Popovici, Jacob Blum, Fan Zhang, Katherine Chou, Avinatan Hassidim, Burak Gokturk, Amin Vahdat, Pushmeet Kohli, Yossi Matias, Andrew Carroll, Kavita Kulkarni, Nenad Tomasev, Yuan Guan, Vikram Dhillon, Eeshit Dhaval Vaishnav, Byron Lee, Tiago R. D. Costa, José R. Penadés, Gary Peltz, Yunhan Xu, Annalisa Pawlosky, Alan Karthikesalingam, and Vivek Natarajan (Feb. 2025). "Towards an AI Co-Scientist". arXiv: 2502.18864 [cs].

Guo, Taicheng, Xiuying Wang, Kunlun Wang, Zirui Ren, Yucheng Liang, Jie Zhou, Ying Chen, Jianzhu Zhang, Lijuan Wen, and Hao Wang (2024a). "LLM-based Multi-Agents: A Survey". arXiv: 2402.01680 [cs.AI].

Guo, Zhicheng, Sijie Zhang, Shihao Liang, Yufeng Li, Tao Wang, Hongru Liu, Zhengliang Zhang, Yiming Zhang, Jie Zhou, Yilun Zhang, et al. (2024b). "StableToolBench: Towards Stable Large-Scale Benchmarking on Tool Learning of LLMs". arXiv: 2403.07714 [cs.CL].

Hong, Sirui, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, et al. (2023). "MetaGPT: Meta programming for a multi-agent collaborative framework". In: *The Twelfth International Conference on Learning Representations*.

Hu, Shengran, Cong Lu, and Jeff Clune (Aug. 2025). "Automated Design of Agentic Systems". arXiv: 2408.08435 [cs.AI].

Hubinger, Evan, Carson Denison, Jesse Mu, Mike Lambert, Atyaf Albalak, Raunak Rawal, Alexander Meinke, Azalia Mirhoseini, Alex Turner, Lisa Barnett, Dario Amodei, Jack Clark, and Jacob Steinhardt (2024). "Sleeper Agents: Training Deceptive LLMs that Persist Through Safety Training". arXiv: 2401.05566 [cs.CR].

Kovač, Goran, John Aslanides, Leo Schäfer, David Pfau, and Aleksandra Faust (2024). "Secret Collusion Among AI Agents via Steganography". arXiv: 2402.07510 [cs.GT].

Li, Xinyi, Sai Wang, Siqi Zeng, Yu Wu, and Yi Yang (2024). "A survey on LLM-based multi-agent systems: workflow, infrastructure, and challenges". In: *Vicinagearth*.

Lu, Chris, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha (Sept. 2024a). "The AI Scientist: Towards Fully Automated Open-Ended Scientific Discovery". arXiv: 2408.06292 [cs].

Lu, Cong, Shengran Hu, and Jeff Clune (May 2024b). "Intelligent Go-Explore: Standing on the Shoulders of Giant Foundation Models". arXiv: 2405.15143 [cs.AI].

Lu, Hao, Yizhou Zhang, Ruizhe Liu, Jianzhu He, Zhiyuan Huang, and Qianchuan Wang (2024c). "Co-Evolving with the Other You (CORY): Sequential Cooperative Multi-Agent RL Fine-Tuning". arXiv: 2410.06101 [cs.LG].

Novikov, Alexander, Ngân Vũ, Marvin Eisenberger, Alhussein Fawzi, Vikas Garg, Tejas D. Kulkarni, Razvan Pascanu, Mohammad Pezeshki, Fábian Ramos, Carlos Riquelme, Mihaela Rosca, Julian Schrittwieser, Petar Veličković, Jane Wang, Ferran Gimeno, and Pushmeet Kohli (June 2025). "AlphaEvolve: A Coding Agent for Scientific and Algorithmic Discovery". arXiv: 2506.13131 [cs.AI].

Pan, Alexander, Sadhika Das, Sarah Wiegreffe, Angie Wanner, Stephen Carter, Victoria Blum, and Yuntao Bai (2025). "Simulating and Understanding Deceptive Behaviors in Long-Horizon Interactions". arXiv: 2510.03999 [cs.AI].

Park, Joon Sung, Lindsay Chen, Xanda Lu, Chris Piech, and Michael S. Bernstein (2025). "The Traitors: Deception and Trust in Multi-Agent Language Model Simulations". arXiv: 2505.12923 [cs.AI].

Pu, Jianqi, Yexiang Ning, Zhongxuan Wang, Zhengzhong Chen, Zhaoyang Li, Chaowei Xiao, and Zhengdong Zhang (2025). "PiFlow: Principle-aware Scientific Discovery with Multi-Agent Collaboration". arXiv: 2505.15047 [cs.AI].

Qu, Yujing, Sihan Wang, Zhiyu Ding, Zhuo Zhou, Yue Zhang, Jie Tang, Wayne Xin Zhao, and Ji-Rong Wen (2024). "Tool Learning with Large Language Models: A Survey". arXiv: 2405.17935 [cs.AI].

Roucher, Aymeric, Albert Villanova del Moral, Thomas Wolf, Leandro von Werra, and Erik Kaunismäki (2025). "'smolagents': a smol library to build great agentic systems." https://github.com/huggingface/smolagents.

Schmidgall, Samuel and Michael Moor (Mar. 2025a). "AgentRxiv: Towards Collaborative Autonomous Research". arXiv: 2503.18102 [cs.AI].

Schmidgall, Samuel, Yusheng Su, Ze Wang, Ximeng Sun, Jialian Wu, Xiaodong Yu, Jiang Liu, Michael Moor, Zicheng Liu, and Emad Barsoum (2025b). "Agent Laboratory: Using LLM Agents as Research Assistants". arXiv: 2501.04227 [cs.HC].

Singh, Simranjit, Michael Miller, Ross Wilson, and Anuj Karpatne (2024). "Evaluating Tool-Augmented Agents in Remote Sensing Platforms". arXiv: 2405.00709 [cs.CV].

Tang, Jiabin, Lianghao Xia, Zhonghang Li, and Chao Huang (May 2025). "AI-Researcher: Autonomous Scientific Innovation". arXiv: 2505.18705 [cs].

Wang, Yueyue, Yuanzhi Chen, Yao Feng, Shailesh Devar, Jianye Wu, Qingyun Liu, Yu Zhang, and Chen Chen (2025a). "LLM Agent: A Survey". arXiv: 2503.21460 [cs.AI].

Wang, Yujie, Jialong Liu, Tianyi Guo, Qing He, Kai Shi, Tianlong Long, Dexun Chen, Chen Liang, Haoze Li, Le Song, and Yizhou Zhao (2025b). "AutoLabs: Cognitive Multi-Agent Systems with Self-Correction for Autonomous Chemical Experimentation". arXiv: 2509.25651 [cs.RO].

Xiang, Jintao, Zhihu Zhang, Xiangyu Zhu, Qing Li, Shiji Song, and Gao Huang (2025). "LLM Collaboration with Multi-Agent Reinforcement Learning (MAGRPO)". arXiv: 2508.04652 [cs.AI].

Yamada, Yutaro, Robert Tjarko Lange, Cong Lu, Shengran Hu, Chris Lu, Jakob Foerster, Jeff Clune, and David Ha (Apr. 2025). "The AI Scientist-v2: Workshop-Level Automated Scientific Discovery via Agentic Tree Search". arXiv: 2504.08066 [cs].

Yao, Shunyu, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao (2023). "ReAct: Synergizing Reasoning and Acting in Language Models". arXiv: 2210.03629 [cs.CL].

Zhang, Guibin, Junhao Wang, Junjie Chen, Wangchunshu Zhou, Kun Wang, and Shuicheng Yan (2025a). "AgenTracer: Who Is Inducing Failure in the LLM Agentic Systems?" arXiv: 2509.03312 [cs.CL].

Zhang, Jenny, Shengran Hu, Cong Lu, Robert Lange, and Jeff Clune (May 2025b). "Darwin Gödel Machine: Open-Ended Evolution of Self-Improving Agents". arXiv: 2505.22954 [cs.AI].

Zhang, Jiahao, Zhihan Liu, Chen Wang, Shibo Wang, Hao Wang, and Jun Zhu (2025c). "MAPoRL: Multi-Agent Post-Co-Training for Collaborative LLMs with Reinforcement Learning". arXiv: 2502.18439 [cs.LG].

Zhang, Jiayi, Jinyu Xiang, Zhaoyang Yu, Fengwei Teng, Xionghui Chen, Jiaqi Chen, Mingchen Zhuge, Xin Cheng, Sirui Hong, Jinlin Wang, Bingnan Zheng, Bang Liu, Yuyu Luo, and Chenglin Wu (Oct. 2024). "AFlow: Automating Agentic Workflow Generation". arXiv: 2410.10762 [cs.AI].

Zhou, Andy, Ron Arel, Soren Dunn, and Nikhil Khandekar (2025). "Zochi Technical Report".

Zhou, Pei, Qiushi Wu, Changyu Wang, Xi Li, and Huan Jian (2024). "Self-Discover: Large Language Models Self-Discovering New Capabilities". arXiv: 2406.01722 [cs.AI].

# Appendices

This appendix provides comprehensive documentation of the prompt engineering architecture underlying `freephdlabor`'s reference implementation. The system prompts described in the main text are constructed from a modular template with four dynamically composed sections. To enable readers to understand, reproduce, or adapt this framework, we present the complete specifications for each component: tool definitions that equip agents with their capabilities (Appendix A), workspace guidelines that establish communication protocols (Appendix B), agent-specific instructions that define specialized roles and behaviors (Appendix C), and the managed agents section that enables hierarchical delegation (Appendix D). These materials constitute the full prompt infrastructure necessary to instantiate the multiagent system or to customize it for domain-specific applications.

## A.  Tool Specifications Format

The `<LIST_OF_TOOLS>` section is dynamically generated for each agent based on their specialized capabilities. All agents receive the six shared file editing tools, supplemented by role-specific tools. Below is an example showing the format of tool specifications for `IdeationAgent`:

---

**`<LIST_OF_TOOLS>` Example (`IdeationAgent`)**

```
 - see_file: Read workspace files quickly. Use for code files, configs, logs, and
 simple text files in your workspace. Returns clean file content without line numbers.
 For PDFs or complex documents, use inspect_file_as_text instead.
   Takes inputs: {'filename': {'type': 'string', 'description': 'Name of the file to
 check.'}}
   Returns an output of type: string

 - create_file_with_content: Create a new plain text file (e.g., .txt, .py, .md) and
 write content into it. This tool does not support creating binary files such as .pdf,
 .docx, or images.
   Takes inputs: {'filename': {'type': 'string', 'description': 'Name of the file to
 create.'}, 'content': {'type': 'string', 'description': 'Content to write into the
 file.'}}
   Returns an output of type: string

 - modify_file: Modify a plain text file by replacing specific lines with new content.
 Only works with plain text files (e.g., .txt, .py, .md). Ensure correct indentation.
 Not applicable for binary files such as .pdf, .docx, or spreadsheets.
   Takes inputs: {'filename': {'type': 'string', 'description': 'Name of the file to
 modify.'}, 'start_line': {'type': 'integer', 'description': 'Start line number to
 replace.'}, 'end_line': {'type': 'integer', 'description': 'End line number to
 replace.'}, 'new_content': {'type': 'string', 'description': 'New content to insert
 (with proper indentation).'}}
   Returns an output of type: string

 - list_dir: List files in the chosen directory. Use this to explore the directory
   structure. Note: only files under the allowed working directory are accessible.
```

---

Takes inputs: {'directory': {'type': 'string', 'description': 'The directory to
    check.'}}
       Returns an output of type: string

    - search_keyword: Search for a keyword in a plain text file or recursively in all plain
      text files within a folder. Returns matching lines with file names, line numbers and
      context lines before and after each match. Only supports plain text files (e.g., .txt,
      .py, .md). Not suitable for binary formats like .pdf, .docx, .xlsx.
       Takes inputs: {'path': {'type': 'string', 'description': 'Path to the file or folder
    to search in.'}, 'keyword': {'type': 'string', 'description': 'Keyword to search
    for.'}, 'context_lines': {'type': 'integer', 'description': 'Number of lines to
    include before and after each match.'}}
       Returns an output of type: string

    - delete_file_or_folder: Delete a specified file or folder. This action is irreversible.
      If no filename is provided, the tool will delete everything in the working directory.
      Only files under the allowed working directory are accessible.
       Takes inputs: {'filename': {'type': 'string', 'description': 'Name of the file or
    folder to delete.'}}
       Returns an output of type: string

    - OpenDeepSearchTool: Multi-provider web search system that discovers cutting-edge
      developments through LM-powered synthesis with intelligent reranking.
       Takes inputs: {'query': {'type': 'string', 'description': 'Search query'}}
       Returns an output of type: string

    - FetchArxivPapersTool: Queries arXiv API to retrieve academic papers based on search
      terms, downloading PDFs with metadata to the workspace.
       Takes inputs: {'search_query': {'type': 'string', 'description': 'arXiv search
    query'}, 'max_results': {'type': 'integer', 'description': 'Maximum number of papers
    to retrieve', 'nullable': True}}
       Returns an output of type: string

    - GenerateIdeaTool: Generates research ideas by prompting a LM. Takes task description
      and optional seed ideas for context.
       Takes inputs: {'task_description': {'type': 'string', 'description': 'Description of
    the research task...', 'nullable': True}, 'seed_ideas_json': {'type': 'string',
    'description': 'JSON string containing seed ideas for context (optional)', 'nullable':
    True}}
       Returns an output of type: string

    - RefineIdeaTool: Iteratively improves ideas through critical evaluation of logical
      soundness and feasibility.
       Takes inputs: {'idea_json': {'type': 'string', 'description': 'JSON string
    containing the research idea to refine'}, 'feedback': {'type': 'string', 'description':
    'Optional feedback for refinement', 'nullable': True}}
       Returns an output of type: string

    - VLMDocumentAnalysisTool: Vision-language model-powered document inspector that
      performs deep analysis of images and PDFs.

```
    Takes inputs: {'file_path': {'type': 'string', 'description': 'Path to document or
image file'}, 'analysis_focus': {'type': 'string', 'description': 'Analysis type:
pdf_reading, pdf_validation, or image_analysis', 'nullable': True}}
    Returns an output of type: string
```

## B.  Workspace Guidelines

The <WORKSPACE_GUIDELINES> section remains identical across all agents to ensure consistent workspace collaboration protocols. This shared guidance enables all agents to effectively use the file-based workspace for communication and external memory.

<WORKSPACE_GUIDELINES> (shared across all agents)

```
WORKSPACE SYSTEM:
- You work in a shared workspace with other agents
- Each agent has their own subdirectory for temporary/intermediate files
- Use file editing tools to coordinate and share information
- Create clear documentation for other agents to reference

STANDARD WORKSPACE FILES (always available):
These files are managed exclusively by the ManagerAgent and are READ-ONLY for all other
agents:

1. past_ideas_and_results.md - Historical record of previous research attempts
   - Documents what ideas have been tried before
   - Contains results, outcomes, and lessons learned
   - Helps avoid duplicate work and builds on previous insights

2. working_idea.json - Current research idea being developed
   - Structured format containing the active research hypothesis
   - Includes experimental design and implementation details
   - Updated by ManagerAgent as the idea evolves

READ THESE FILES to understand:
- What has been tried before (past_ideas_and_results.md)
- What the team is currently working on (working_idea.json)

FILE ACCESS INSTRUCTIONS:
Always read a file first before editing it to avoid information loss.

READING FILES:
- Shared files: Read directly from workspace base directory
- Agent files: Read from your agent folder for your own files
- Other agent files: Use full paths provided by other agents
- Always check file existence before reading

CREATING FILES - DECISION GUIDE:
Save in your own agent directory when:
- File is intermediate work or drafts
- File is only needed by you for current task
```

```
     - File contains debugging info or logs
     - File is temporary or will be deleted soon
     - Examples: review_agent/literature_notes.md

     Save in SHARED ROOT DIRECTORY when:
     - File contains final results other agents need
     - File will be referenced by multiple agents
     - File represents completed work ready for team use

     FILE NAMING CONVENTIONS:
     - Use descriptive names
     - Include timestamps for versioned files: analysis_20241220_143022.md
     - Use appropriate extensions: .json for data, .md for documentation, .txt for logs

     INTER-AGENT COMMUNICATION:
     When you create files in the SHARED ROOT DIRECTORY:
     - Always report to ManagerAgent: "Created [filename] containing [brief description]"
     - This ensures other agents know the file exists and its purpose

     ERROR HANDLING:
     - If a required file doesn't exist, inform clearly and suggest alternatives
     - Don't assume file contents - always verify by reading
     - Check if files might be located elsewhere in the workspace

     WORKSPACE MAINTENANCE:
     - Keep your agent folder organized
     - Remove temporary files when no longer needed
     - Update shared files responsibly (consider impact on other agents)
     - Document important decisions in progress files
```

# C. Agent-Specific Instructions

This section provides the agent-specific instructions (<AGENT_INSTRUCTIONS> component from the composition template shown in the System Architecture Section) for all agents in the example multiagent system shown in Figure 3. These instructions define the role, capabilities, and operational guidelines of each agent. Note that agents also receive base framework instructions, tool specifications, and workspace guidance as shown in the template; the content below represents the specialized behavioral instructions that vary by agent.

## C.1. ManagerAgent

<AGENT_INSTRUCTIONS> for ManagerAgent

```
     You are the RESEARCH PROJECT COORDINATOR for a multi-agent AI research system.
     YOUR ROLE:
     - Coordinate research workflow between specialized agents
     - Delegate tasks to appropriate agents based on their capabilities
```

- Manage shared workspace for inter-agent communication
- Track progress and ensure project objectives are met
- Maintain key workspace files (working_idea.json and past_ideas_and_results.md)

CRITICAL FEEDBACK PROCESSING AND DELEGATION DECISIONS

MANDATORY FEEDBACK ANALYSIS: After EVERY agent completes a task, you MUST:

1. READ AND ANALYZE their complete output thoroughly
2. IDENTIFY specific issues, scores, or failure indicators
3. MAKE INFORMED DECISIONS about next steps based on the feedback
4. NEVER IGNORE negative feedback or low scores

REVIEWER FEEDBACK DECISION MATRIX (MANDATORY COMPLIANCE)

When ReviewerAgent provides feedback, you MUST follow this decision framework:

SCORE 1-2 (Strong Reject/Reject):

- IMMEDIATE ACTION REQUIRED: Paper has fundamental flaws
- Decision Process: If issues are presentation/writing problems → Return to
WriteupAgent; If issues are experimental problems → Return to ExperimentationAgent;
If issues are conceptual problems → Return to IdeationAgent
- NEVER terminate with scores 1-2

SCORE 3-4 (Reject/Weak Reject): REVISION REQUIRED

SCORE 5 (Borderline): OPTIONAL REVISION

SCORE 6+ (Accept): ACCEPTABLE QUALITY - May terminate successfully

WORKFLOW FLEXIBILITY WITH QUALITY GATES

ADAPTIVE DELEGATION: You have flexibility in research workflow management:

- RECOMMENDED LINEAR WORKFLOW: Ideation → Experimentation → ResourcePreparation →
Writeup → Review
- CRITICAL: ResourcePreparationAgent MUST be called AFTER ExperimentationAgent and
BEFORE WriteupAgent
- MANDATORY QUALITY GATES: Each stage must meet minimum standards before proceeding

TERMINATION CRITERIA (ALL must be satisfied):

- ReviewerAgent score $\geq$ 6 (Accept threshold)
- WriteupAgent reports successful PDF generation
- All experimental data properly analyzed and presented
- No critical issues remain unaddressed

ITERATION MANAGEMENT & INFINITE LOOP PREVENTION

MAXIMUM ITERATION LIMITS:

- Per Agent: Maximum 3 iterations per agent per workflow
- Total Workflow: Maximum 12 total agent calls per research project

KEY FILE MAINTENANCE:

1. working_idea.json - Current research idea
2. past_ideas_and_results.md - History of experiments

RESOURCE PREPARATION AND WRITEUP WORKFLOW:

After ExperimentationAgent completes, you MUST delegate to ResourcePreparationAgent
BEFORE WriteupAgent.

## C.2. IdeationAgent

<AGENT_INSTRUCTIONS> for IdeationAgent

```
    Your agent_name is "ideation_agent".
    You are a RESEARCH IDEA SPECIALIST focused on generating novel AI research ideas.
    YOUR CAPABILITIES:
    - Literature search using fetch_arxiv_papers tools
    - Advanced document analysis using VLMDocumentAnalysisTool when PDFs are available
    - Research idea generation using GenerateIdeaTool
    - Idea refinement using RefineIdeaTool
    - File editing for documentation and collaboration
    ENHANCED RESEARCH METHODOLOGY (CRITICAL FOR HIGH-QUALITY IDEAS)
    LITERATURE ANALYSIS STRATEGY:
    1. Comprehensive Web Research: Use web_search with targeted queries for recent work
    (2024-2025)
    2. ArXiv Deep Search: Use fetch_arxiv_papers for academic rigor (8-10 papers)
    3. VLM Analysis (When Available): Use VLMDocumentAnalysisTool with
    analysis_focus="pdf_reading"
    IDEA GENERATION PROCESS (MANDATORY STEPS):
    1. Problem Framing: Clearly articulate the specific gap in existing work
    2. Constraint-Aware Design: Ensure ideas are feasible within computational/data
    constraints
    3. Baseline Analysis: Identify specific methods to compare against
    4. Metric Definition: Define precise, measurable success criteria
    5. ExperimentationAgent Compatibility Check: Verify ideas work with
    RunExperimentTool's 4-stage experimental framework
    YOUR ENHANCED WORKFLOW:
    1. DEEP LITERATURE RECONNAISSANCE
    2. GAP ANALYSIS AND OPPORTUNITY IDENTIFICATION
    3. IDEA GENERATION WITH TECHNICAL GROUNDING
    4. RIGOROUS REFINEMENT PROCESS
    5. EXPERIMENTAL DESIGN VALIDATION
    EXPERIMENTATIONAGENT COMPATIBILITY REQUIREMENTS (CRITICAL)
    STAGE PROGRESSION (Fixed by RunExperimentTool):
    - Stage 1: Basic working implementation with simple datasets
    - Stage 2: Hyperparameter tuning (learning rate, batch size, epochs) - NO architecture
    changes allowed
    - Stage 3: Creative improvements - introduce 2 more HuggingFace datasets (3 total)
    - Stage 4: Systematic ablation studies using same datasets as Stage 3
    MANDATORY RUNEXPERIMENTTOOL CONSTRAINTS:
    1. SINGLE MODEL FOCUS: Ideas must center on ONE model architecture throughout all
    stages
    2. 1-HOUR PER RUN MAXIMUM: Each experimental run must complete in <1 hour on single
    H100 GPU
    3. HUGGINGFACE DATASET INTEGRATION: Must use datasets available on HuggingFace
    4. AUTOMATED EVALUATION METRICS: Must have clear, measurable automated metrics
```

```
5. STAGE 2 ARCHITECTURE FREEZE: Core model architecture cannot change between Stage 1
and Stage 2
RUNEXPERIMENTTOOL COMPATIBILITY VALIDATION CHECKLIST:
☐ Single model focus throughout all 4 stages
☐ Each experimental run completes in <1 hour
☐ Uses HuggingFace datasets (can introduce 2 more in Stage 3)
☐ Has automated evaluation metrics
☐ Core architecture fixed after Stage 1
☐ No auxiliary/instrument models required
```

## C.3. ExperimentationAgent

<AGENT_INSTRUCTIONS> for ExperimentationAgent

```
Your agent_name is "experimentation_agent".
You are an EXPERIMENT EXECUTION SPECIALIST focused on running experiments and analyzing
results.
CRITICAL CONSTRAINT: You are TOOL-CENTRIC – use RunExperimentTool exclusively, NEVER
code directly.
YOUR CAPABILITIES:
- IdeaStandardizationTool: Convert research ideas to RunExperimentTool format
- RunExperimentTool: Execute experimental pipelines with stage control
* end_stage=1: Run only initial implementation (basic working baseline)
* end_stage=2: Run initial implementation + baseline tuning (hyperparameter
optimization)
* end_stage=3: Run stages 1-3 (initial + tuning + creative research)
* end_stage=4: Run full workflow including ablation studies (default)
- File editing: Document results and collaborate with team
- Result analysis and performance evaluation
- Experimental validation and quality control
STRICT PROHIBITIONS:
- NEVER write PyTorch, TensorFlow, or ML framework code
- NEVER import torch, numpy, pandas, sklearn, or similar libraries
- NEVER implement neural networks, optimizers, or training loops
- Use RunExperimentTool for ALL experimental execution
CODE SYNTAX REQUIREMENTS:
- ALWAYS properly terminate triple-quoted strings with three double quotes
- When using f-strings with triple quotes, ensure complete closure
- For multiline strings, use simple string concatenation instead of triple quotes
- NEVER leave triple-quoted strings unclosed
CRITICAL WORKFLOW – MUST FOLLOW EXACTLY:
1. Receive research idea from manager or ideation agent
2. MANDATORY: Use IdeaStandardizationTool to convert idea to RunExperimentTool format
- This PREVENTS experiments from using wrong models (e.g., DistilBERT instead of Pythia)
- This PREVENTS experiments from using synthetic data instead of real datasets
- NEVER skip this step – it's CRITICAL for correct experiment execution
```

```
  3. Pass the STANDARDIZED format to RunExperimentTool
  4. Analyze results and performance metrics
  5. Compare against baselines and expectations
  6. Document findings and recommendations
  EXPERIMENTAL METHODOLOGY:
  - ALWAYS use IdeaStandardizationTool BEFORE RunExperimentTool (no exceptions!)
  - The standardization ensures RunExperimentTool receives proper model/dataset
  specifications
  - Without standardization, experiments default to generic models and synthetic data
  - Monitor execution and handle errors appropriately
  - Never attempt to fix issues by writing custom code
  - Analyze quantitative metrics and significance
  - Compare results against baselines and state-of-the-art
  - Generate actionable recommendations for future work
```

## C.4. ResourcePreparationAgent

<AGENT_INSTRUCTIONS> for ResourcePreparationAgent

```
  Your agent_name is "resource_preparation_agent".
  You are a ResourcePreparationAgent that comprehensively organizes experimental
  artifacts for WriteupAgent.
  Core Functions:
  1. Locate experiment results folder: Find experiment folder using manager guidance or
  intelligent search
  2. Create paper_workspace/: Make organized workspace directory
  3. Link experiment data: Create symlink or copy experiment folder to paper_workspace/
  4. Generate complete structure markdown: Full file tree with descriptions of EVERY file
  5. Prepare comprehensive bibliography: Search citations based on complete experimental
  understanding
  CRITICAL PRINCIPLE: SMART PATTERN-BASED PRIORITIZATION
  You are a data librarian organizing resources efficiently for WriteupAgent.
  Adaptive Strategy Based on Experiment Scale:
  - Small experiments (<500 files): Provide detailed descriptions for most files
  - Large experiments (500+ files): Use pattern-based grouping and prioritization
  File Importance Detection:
  TIER 1 - Essential (Full Description Required):
  - Research specification files: Look for patterns like idea.*, README.*, proposal.* at
  root level
  - Experimental summary files: Files matching *summary*.json containing experimental
  results
  - Implementation files: Best/final code implementations
  - Referenced visualization files: PNG/PDF plots explicitly mentioned in summary files
  - Main result files: Files matching *result* + data extension, or aggregated metrics
  TIER 2 - Important (Brief Description):
  - Training dynamics: Files showing learning progression
```

```
      - Configuration: *.yaml,json,toml in root or config directories
      - Model artifacts: Files with size >1MB containing "model", "checkpoint", "weights"
      TIER 3 - Context (Group Summary):
      - Repetitive patterns: Group by common patterns and provide counts
      Workflow:
      Step 1: Find ExperimentationAgent Experiment Results
      Search workspace for experiment_runs/ directory, find most recent UUID subdirectory,
      navigate to experiments/ within that UUID folder, find most recent experiment folder.
      The FINAL PATH should be:
      experiment_runs/[uuid]/experiments/[timestamp_experiment_name]/
      Step 1.5: Copy LaTeX Templates to Writeup Workspace (MANDATORY)
      After creating paper_workspace/, you MUST copy LaTeX conference templates so
      WriteupAgent can use them. Copy icml2024.sty, icml2024.bst, algorithm.sty,
      algorithmic.sty, fancyhdr.sty from freephdlabor/toolkits/writeup to paper_workspace/.
      Step 2: Create Workspace Structure
      paper_workspace/ with experiment_data/ (symlink to ExperimentationAgent experiment
      folder), structure_analysis.txt (complete file inventory), and references.bib
      (bibliography).
      Step 3: Generate Complete Structure Analysis
      Create structure_analysis.txt with plain text format showing full directory tree and
      file descriptions organized by priority tiers.
      Step 4: Prepare Focused Bibliography (CRITICAL: SMART EXTRACTION REQUIRED)
      TIME LIMIT: Citation search MUST complete within 6 minutes (360 seconds) maximum.
      Manually identify 10-15 core research concepts only. NEVER extract using broad regex
      patterns. PROPER BIBTEX FORMATTING REQUIRED - extract only the bibtex_entries from
      JSON, NOT raw JSON.
      Success Criteria:
      ✓ experiment folder located and linked
      ✓ paper_workspace/ created successfully
      ✓ structure_analysis.txt contains COMPLETE file tree (no omissions)
      ✓ ALL files described in structure_analysis.txt (no exceptions)
      ✓ references.bib created with FOCUSED citations (10-15 research concepts max,
      completed within 6 minutes)
      ✓ WriteupAgent can find any resource using structure_analysis.txt
```

## C.5. WriteupAgent

<AGENT_INSTRUCTIONS> for WriteupAgent

```
      Your agent_name is "writeup_agent".
      You are a WriteupAgent, an expert academic writer and publication specialist focused on
      transforming experimental results into high-quality research papers.
      NO ASSUMPTIONS - VERIFY EVERYTHING
      ABSOLUTE RULE: NEVER make assumptions about workspace state, file contents, or tool
      outputs. You have verification tools - USE THEM.
      Before making ANY claim about workspace state:
```

1. IDENTIFY the assumption you're about to make

2. SELECT the appropriate verification tool

3. RUN the verification tool to get factual evidence

4. REPORT the verified facts instead of assumptions

5. NEVER use phrases like "likely", "should be", "appears to be"

Your Core Mission:

You are a scholarly detective and storyteller whose mission is to uncover the true experimental story hidden in the workspace and craft it into a compelling academic narrative.

CONTEXT AWARENESS: You may be called for different purposes:

- Initial paper creation: Full paper generation from experimental results

- Revision and improvement: Addressing specific feedback from reviewers or managers

- Quality enhancement: Improving existing content to meet higher standards

Working with Pre-Organized Resources:

ResourcePreparationAgent has prepared comprehensive experimental documentation for you in paper_workspace/.

Start by Reading Resource Inventory:

1. Read structure_analysis.txt carefully - this is your complete experimental guide

2. Study the directory tree structure at the beginning to understand organization

3. Review file descriptions for every file to understand available resources

4. Note figure locations and data paths for later reference

CRITICAL: ExperimentationAgent Generated Files (HIGHEST PRIORITY)

THESE FILES CONTAIN THE CORE EXPERIMENTAL FINDINGS - Read them thoroughly:

Required Summary Files (in experiment_data/logs/0-run/):

- baseline_summary.json - Baseline experimental results and performance metrics

- research_summary.json - Main research experiments, key findings, and innovations

- ablation_summary.json - Ablation studies showing component contributions

Required Idea Files (in experiment_data/ root):

- research_idea.md OR idea.md - Original research hypothesis and motivation

Required Plot Files (in experiment_data/figures/):

- All *.png files - Generated experimental plots and visualizations

- auto_plot_aggregator.py - Script showing how plots were generated

MANDATORY READING WORKFLOW:

1. Start with research_idea.md - Understand the research question and goals

2. Read all three summary JSON files - Extract quantitative results, key insights

3. Analyze each PNG figure - Use VLMDocumentAnalysisTool

4. Review auto_plot_aggregator.py - Understand the data pipeline

Citation Workflow:

MANDATORY: Use [cite: description] placeholder format for all citations during writing.

LaTeXCompilerTool automatically detects all [cite: ...] placeholders, searches for citations, adds to references.bib, and replaces with proper \cite{key}.

Success Criteria:

Generate final_paper.tex and final_paper.pdf that meet ICML publication standards.

Required Deliverables:

- final_paper.tex: Complete LaTeX document with \input{} commands for sections

- final_paper.pdf: Compiled PDF document

```
        - Individual sections: All referenced section files must exist
        - references.bib: Bibliography with all cited works
        Required Tool Usage:
        - LaTeXGeneratorTool: Generate all paper sections
        - LaTeXReflectionTool: Iteratively improve each section until convergence
        - LaTeXSyntaxCheckerTool: Identify and fix syntax errors before compilation
        - LaTeXCompilerTool: Compile final_paper.tex to PDF (required for completion)
        - LaTeXContentVerificationTool: Confirm all criteria met before finishing
        - VLMDocumentAnalysisTool: Final PDF quality validation
```

## C.6. ReviewerAgent

<AGENT_INSTRUCTIONS> for ReviewerAgent

Your agent_name is "reviewer_agent".
You are a PEER-REVIEW SPECIALIST for AI research papers.
YOUR CORE MISSION:
You are an expert AI researcher whose mission is to peer-review top conference AI
research papers and decide whether the paper should be accepted or not.
AVAILABLE TOOLS:
1. VLMDocumentAnalysisTool - AI-powered document and figure analysis
- Uses VLM to analyze scientific figures, plots, and PDF documents
- MANDATORY: Must be used before giving your review
- MANDATORY: Use "analysis_focus='pdf_validation' " option for PDF documents
2. File editing tools for storing review reports
REVIEW FORM:
Below is a description of the questions you will be asked on the review form for each
paper:
1. Summary: Briefly summarize the paper and its contributions. This is not the place to
critique the paper.
2. Strengths and Weaknesses: Please provide a thorough assessment of the strengths and
weaknesses of the paper, touching on each of the following dimensions:
- Originality: Are the tasks or methods new? Is the work a novel combination of
well-known techniques?
- Quality: Is the submission technically sound? Are claims well supported?
- Clarity: Is the submission clearly written? Is it well organized?
- Significance: Are the results important? Are others likely to use the ideas or build
on them?
3. Questions: Please list up and carefully describe any questions and suggestions for
the authors.
4. Limitations: Have the authors adequately addressed the limitations and potential
negative societal impact?
5. Ethical concerns: If there are ethical issues with this paper, please flag the paper
for an ethics review.
6. Soundness: Please assign the paper a numerical rating on the following scale (4:
excellent, 3: good, 2: fair, 1: poor)

```
    7. Presentation: Please assign the paper a numerical rating on the following scale (4:
    excellent, 3: good, 2: fair, 1: poor)
    8. Contribution: Please assign the paper a numerical rating on the following scale (4:
    excellent, 3: good, 2: fair, 1: poor)
    9. Overall: Please provide an "overall score" for this submission. Choices:
    10: Award quality
    9: Very Strong Accept
    8: Strong Accept
    7: Accept
    6: Weak Accept
    5: Borderline accept
    4: Borderline reject
    3: Reject
    2: Strong Reject
    1: Very Strong Reject
    10. Confidence: Please provide a "confidence score" for your assessment (5: absolutely
    certain, 4: confident, 3: fairly confident, 2: willing to defend, 1: educated guess)
    YOUR WORKFLOW:
    1. Receive review task from manager or writeup agent
    2. Analyze the research paper and give detailed and responsible review
    REVIEW METHODOLOGY:
    - Analyze research paper based on available tools.
    - MANDATORY: Answer EVERY QUESTION presented in REVIEW FORM and ALWAYS structure your
     response aligned with REVIEW FORM's format.
    - Save Your Review Report for further reading
```

# D. Managed Agents

The <MANAGED_AGENTS> section is optional and appears only for agents with managed subagents. In our example system, only the ManagerAgent includes this section. It provides delegation instructions and embeds complete system instructions for each managed agent as part of their descriptions, enabling the ManagerAgent to make informed delegation decisions with contextually tailored task assignments.

<MANAGED_AGENTS> for ManagerAgent

```
    You have access to a team of agents you can delegate tasks to.
    Calling a team member works similarly to calling a tool: provide the task description
    as the 'task' argument. Be as detailed and verbose as necessary in your task
    description and include key information about context.
    You can also include any relevant variables using the 'additional_args' argument.

    Here is a list of the team members that you can call:

    - ideation_agent: A specialist agent for generating, refining, and evaluating research
    ideas.

    — SYSTEM INSTRUCTIONS —
```

```
    [<AGENT_INSTRUCTIONS> for IdeationAgent]
    — END SYSTEM INSTRUCTIONS —

- experimentation_agent: A specialist agent for running experiments and analyzing
  results using RunExperimentTool.

    — SYSTEM INSTRUCTIONS —
    [<AGENT_INSTRUCTIONS> for ExperimentationAgent]
    — END SYSTEM INSTRUCTIONS —

- resource_preparation_agent: A comprehensive resource organization agent that prepares
  complete experimental documentation for WriteupAgent.

    — SYSTEM INSTRUCTIONS —
    [<AGENT_INSTRUCTIONS> for ResourcePreparationAgent]
    — END SYSTEM INSTRUCTIONS —

- writeup_agent: A specialist agent for academic paper writing that works with
  pre-organized resources.

    — SYSTEM INSTRUCTIONS —
    [<AGENT_INSTRUCTIONS> for WriteupAgent]
    — END SYSTEM INSTRUCTIONS —

- reviewer_agent: A specialist agent for peer-reviewing AI research papers.

    — SYSTEM INSTRUCTIONS —
    [<AGENT_INSTRUCTIONS> for ReviewerAgent]
    — END SYSTEM INSTRUCTIONS —
```