# Towards Scientific Intelligence: A Survey of LLM-based Scientific Agents

**Shuo Ren**[*1,2], **Pu Jian**[*1,3], **Zhenjiang Ren**[*1,3], **Chunlin Leng**[*1,3],
**Can Xie**[*1,3], **Jiajun Zhang**[†1,2,3,4]

[1] State Key Laboratory of Multimodal Artificial Intelligence Systems,
[2] Foundation Model Research Center, Institute of Automation, CAS.
[3] University of Chinese Academy of Science, Beijing, China.
[4] Wuhan AI Research, Wuhan, China. [†]jjzhang@nlpr.ia.ac.cn
[*]{shuo.ren, jianpu2023, renzhenjiang2024, lengchunlin2023, xiecan2024}@ia.ac.cn

## Abstract

As scientific research becomes increasingly complex, innovative tools are needed to manage vast data, facilitate interdisciplinary collaboration, and accelerate discovery. Large language models (LLMs) are now evolving into LLM-based scientific agents that automate critical tasks—ranging from hypothesis generation and experiment design to data analysis and simulation. Unlike general-purpose LLMs, these specialized agents integrate domain-specific knowledge, advanced tool sets, and robust validation mechanisms, enabling them to handle complex data types, ensure reproducibility, and drive scientific breakthroughs. This survey provides a focused review of the architectures, design, benchmarks, applications, and ethical considerations surrounding LLM-based scientific agents. We highlight why they differ from general agents and the ways in which they advance research across various scientific fields. By examining their development and challenges, this survey offers a comprehensive roadmap for researchers and practitioners to harness these agents for more efficient, reliable, and ethically sound scientific discovery.

## 1 Introduction

Imagine an AI agent that autonomously designs a groundbreaking vaccine, optimizes chemical reactions with pinpoint precision, or uncovers hidden patterns in astronomical data—all while adhering to ethical standards and reproducibility. This is no longer science fiction. Large language models (LLMs), once confined to text generation, are now at the forefront of transforming scientific research by serving as specialized scientific agents that automate complex research tasks such as hypothesis generation, experiment design, and data analysis.

Modern scientific research is becoming increasingly complex, demanding innovative tools that not only manage vast amounts of information but also facilitate interdisciplinary discovery. In response, LLM-based scientific agents have evolved into systems specifically designed for the scientific domain. Unlike general-purpose LLMs, these agents integrate domain-specific knowledge, interface with tailored tools, and process diverse data types—including numerical datasets, chemical structures, and biological sequences. Consequently, they are uniquely positioned to streamline critical research tasks and drive rapid scientific breakthroughs.

As the adoption of these agents grows, a systematic review of their development, applications, and challenges becomes essential. While existing surveys provide comprehensive overviews of general LLM-based agents (Wang et al., 2024b; Xi et al., 2023; Guo et al., 2024; Hu et al., 2024a; Li et al., 2024d; Xie et al., 2024; Cheng et al., 2024; Shen, 2024), focusing specifically on LLM-based scientific agents is crucial given their distinctive roles and requirements in the scientific domain. Here's why this specialized survey is valuable:

**1. Domain-Specific Applications**: Scientific agents are designed specifically for research tasks such as experimental design, data analysis, and hypothesis generation. They incorporate deep scientific methodologies and domain-specific expertise, enabling them to handle the rigorous demands of research workflows—capabilities that general-purpose LLM agents, with their broad and non-specialized approaches, do not possess.

**2. Integration with Scientific Tools**: Unlike general-purpose agents, scientific agents are architected to integrate seamlessly with specialized scientific tools, laboratory instruments, and advanced simulators. This integration supports real-time simulation, precise control, and robust validation of experimental processes, ensuring the agent can manage complex scientific operations.

**3. Handling Complex Scientific Data**: Sci-

---

[*]These authors contribute equally to this work.
[†]Corresponding author.

entific research involves complex data types, including numerical data, chemical structures, and biological sequences. LLM-based scientific agents must be equipped to process and interpret these data forms accurately, a requirement less prevalent in general-purpose LLM agents.

**4. Ethical and Reproducibility Concerns**: Scientific agents adhere to strict ethical standards and incorporate rigorous validation and error-checking mechanisms, such as self-review and statistical analyses, to ensure that their outputs are reliable and reproducible—features typically not addressed by general-purpose LLM agents.

**5. Advancement of Scientific Discovery**: The ultimate goal of scientific agents is to accelerate scientific discovery and innovation. This objective requires capabilities beyond those of general LLM agents, including the ability to generate novel hypotheses, design experiments, and interpret complex results within specific scientific contexts.

By focusing on these distinct aspects, a survey dedicated to LLM-based scientific agents can provide deeper insights into their development, applications, and the unique challenges they face, offering valuable guidance for researchers and practitioners in this specialized field. We hope this survey provides a roadmap for researchers and practitioners to harness these agents effectively, paving the way for faster, more reproducible, and ethically sound scientific discovery.

The remainder of this survey is organized as follows: In **Section 2 (Architectures)**, we begin by examining the fundamental design of these agents. This section is subdivided into three main parts: first, the role of the Planner in decomposing and managing scientific tasks; second, the various Memory mechanisms that enable context retention and iterative learning; and third, the integration of specialized Tool Sets that extend scientific capabilities. After that, in **Section 3 (General-purpose vs. Scientific Agents)**, we will give a detailed comparison between general-purpose and scientific agents, elaborating the reasons why scientific agents need careful design. In **Section 4 (Benchmarks)**, we review the evaluation frameworks used to assess both the general reasoning ability and the scientific research-oriented performance of LLM-based agents. In **Section 5 (Applications)**, we explore real-world applications of LLM-based agents in scientific research, highlighting how these systems are deployed to solve complex problems across various disciplines. In **Section 6 (Ethics)**, we address

the ethical implications and reproducibility challenges inherent in deploying these agents, ensuring that their outputs are not only efficient but also responsible and transparent.

Additionally, we conclude each subsection with a discussion on the challenges and potential directions for future research, offering guidance for both scholars and practitioners in harnessing the full potential of LLM-based scientific agents.
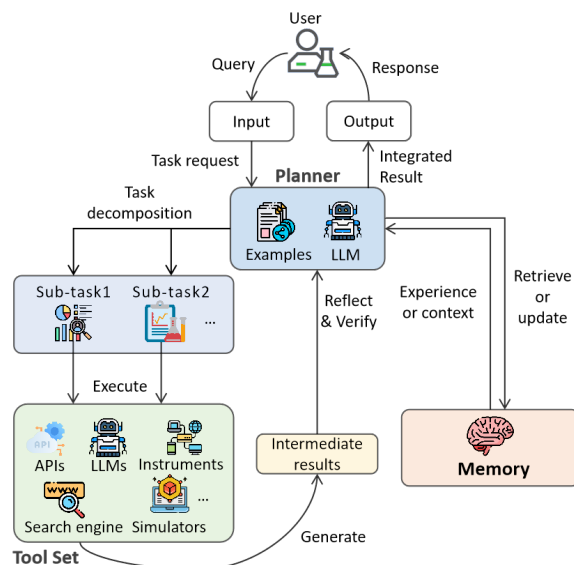


Figure 1: A typical architecture of LLM-based scientific agents. Note that in mainstream agent frameworks, planners are predominantly implemented based on LLMs, and their capabilities include task planning, reflection, and verification, etc. For the sake of abstraction, we represent these functions with a single planner in this architecture diagram. However, in specific implementations, different agents might be set up to accomplish distinct functions (see Section 2.1.6 for further discussion about single-agent planners vs. multi-agent planners).

## 2 Architectures

The architecture of LLM-based scientific agents is designed to enable iterative, context-aware processing of complex scientific tasks. It typically consists of three core components: Planner, Memory, and Tool Set as shown in Figure 1. The workflow begins with the user submitting a query, which is typically a scientific task in the form of text and scientific data. The query is received as input by the system. The Planner decomposes the task into sub-tasks, retrieves relevant context or knowledge from Memory, and executes actions via the Tool Set (e.g., APIs, simulators, instruments, search engines, etc). Note that the LLM itself can also be treated as

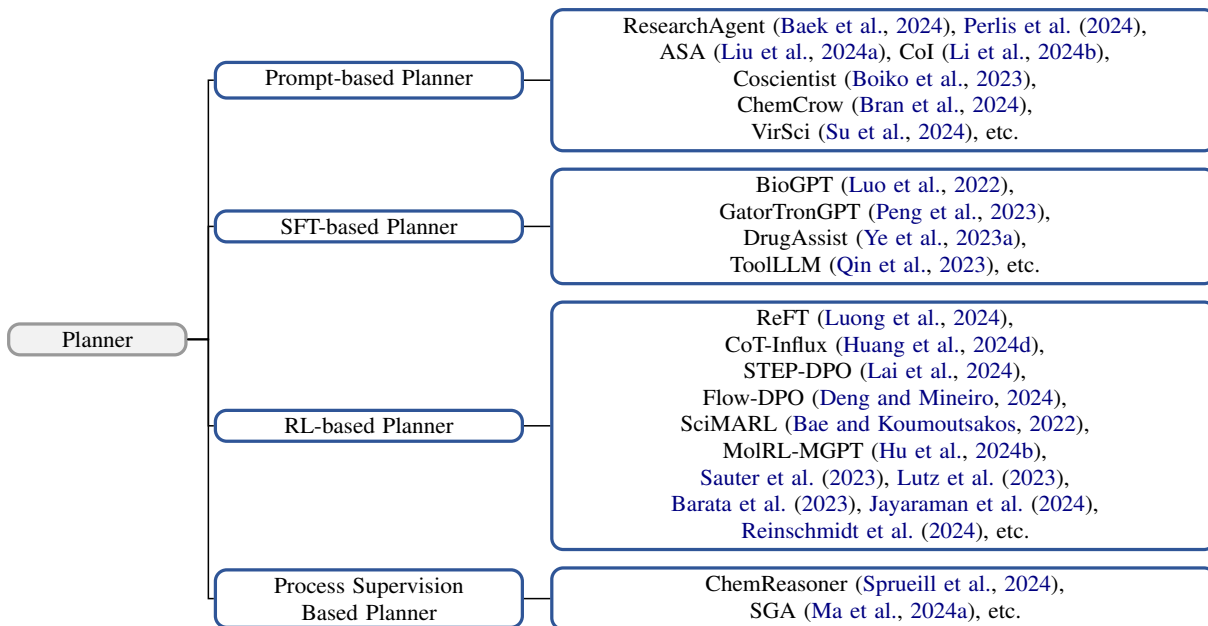| | |
|---|---|
| **Prompt-based Planner** | ResearchAgent (Baek et al., 2024), Perlis et al. (2024), ASA (Liu et al., 2024a), CoI (Li et al., 2024b), Coscientist (Boiko et al., 2023), ChemCrow (Bran et al., 2024), VirSci (Su et al., 2024), etc. |
| **SFT-based Planner** | BioGPT (Luo et al., 2022), GatorTronGPT (Peng et al., 2023), DrugAssist (Ye et al., 2023a), ToolLLM (Qin et al., 2023), etc. |
| **RL-based Planner** | ReFT (Luong et al., 2024), CoT-Influx (Huang et al., 2024d), STEP-DPO (Lai et al., 2024), Flow-DPO (Deng and Mineiro, 2024), SciMARL (Bae and Koumoutsakos, 2022), MolRL-MGPT (Hu et al., 2024b), Sauter et al. (2023), Lutz et al. (2023), Barata et al. (2023), Jayaraman et al. (2024), Reinschmidt et al. (2024), etc. |
| **Process Supervision Based Planner** | ChemReasoner (Sprueill et al., 2024), SGA (Ma et al., 2024a), etc. |

Figure 2: Taxonomy of the planner of science agents.

a tool to finish the related sub-tasks such as reasoning. The actions will generate some intermediate results, which are reflected and verified by the Planner, and Memory is updated with new knowledge to refine future decisions. If the reflection indicates further actions, the Planner will make new plans to do modification. This iterative process continues until the verification passes and the final integrated result is generated, then it is returned to the user as output. Note that while previous LLM-based multi-modal agents often include a separate perceptron module to handle multi-modal inputs (Xie et al., 2024), our survey integrates multi-modal scientific data perception as a fundamental capability of the Planner for the sake of simplification. In the following subsections, we will introduce the three components respectively.

## 2.1 Planner

The Planner serves as the logical core of LLM-based scientific agents, orchestrating structured, method-driven workflows rather than merely decomposing tasks in an ad hoc manner. By integrating domain knowledge with specialized reasoning strategies, the Planner translates high-level scientific problems into reproducible sub-tasks and coordinates their execution across the system. In this context, it enforces hierarchical planning that mirrors the scientific method—defining hypotheses, selecting tools for experimentation or simulation, and validating outcomes before moving forward. Planner designs in scientific agents can be broadly classified into four approaches as in Table 1, i,e, prompt-based, supervised fine-tuning (SFT), reinforcement learning (RL), and process supervision—each offering distinct mechanisms for incorporating domain-specific constraints and robust validation into the research process. The taxonomy of related work is provided in Figure 2.

### 2.1.1 Prompt-Based Planner

This subsection focuses on prompt-based planning of scientific agents, which harness the power of carefully engineered prompts to trigger in-context learning (ICL), thereby guiding the agent to produce a logical, structured and step-by-step plan without requiring additional fine-tuning, as illustrated in Figure 3(a).

Several studies have demonstrated the potential of prompt-based planning in scientific contexts, showcasing the ability of LLM-based scientific agents to tackle complex tasks. We have classified these works into three categories based on different prompt constructions, as shown in Table 2. Note that some works use more than one type of prompt, and we only exemplify their typical type.

**Contextual Information Embedding** is a key approach in prompt-based planning, where detailed, context-specific information is embedded within the prompt to guide the agent's decision-making. For instance, Perlis et al. (2024) create a clinical decision support agent for bipolar disorder, using patient history, symptoms, and guidelines in the

| Approach | Methodology | Strengths and Limitations | Typical Use Cases |
|---|---|---|---|
| Prompt-Based | Use carefully engineered prompts for in-context learning | ✓ No additional training required<br>✓ Flexible and easy to implement<br>✗ Highly dependent on prompt quality<br>✗ May lack robustness for complex tasks | Rapid prototyping; initial task decomposition |
| SFT-Based (Supervised Fine-Tuning) | Fine-tune a pre-trained LLM on curated planning trajectories | ✓ Adapts to domain-specific nuances<br>✓ Produces more precise, step-by-step plans<br>✗ Requires large, high-quality labeled datasets<br>✗ Computationally resource intensive | Detailed planning tasks; scientific workflows with structured data |
| RL-Based (Reinforcement Learning) | Optimize decision-making through reward and penalty signals | ✓ Learns adaptive strategies through trial and error<br>✓ Improves planning over iterations<br>✗ Needs well-defined reward functions<br>✗ Computationally expensive | Iterative, multi-step processes (e.g., experimental design) |
| Process Supervision-Based | Incorporate iterative self-evaluation and external feedback loops | ✓ Mimics human error-correction<br>✓ Continuously refines the planning process<br>✗ More complex to implement<br>✗ May require additional verification mechanisms | Dynamic tasks requiring ongoing refinement and reliability |

Table 1: Comparison of different planners.

| Prompt Type | Method | Prompt Contents | Task |
|---|---|---|---|
| Contextual Information Embedding | Perlis et al. (2024) | Patient's clinical history, symptoms, and guidelines | Clinical decision support |
| | CoI (Li et al., 2024b) | Related research papers with key findings | Research idea generation |
| | ResearchAgent (Baek et al., 2024) | Initial research ideas and background knowledge | Iterative research agent |
| Iterative Feedback and Refinement | ResearchAgent (Baek et al., 2024) | Instructions for review agents to provide feedback | Iterative research agent |
| | LogicSolver (Yang et al., 2022) | Mathematical problems with solution and reasoning instructions | Mathematical problem solving |
| Task Structuring and Multi-Agent Coordination | Coscientist (Boiko et al., 2023) | Four commands to define the action space | Autonomous experimental design and execution |
| | ASA (Liu et al., 2024a) | Instructions for experimental design, simulation, analysis, and report | Automated simulation |
| | Virci (Su et al., 2024) | Role and task descriptions | Generate, evaluate, and refine research idea |
| | ChemCrow (Bran et al., 2024) | Specific instructions about the task and the desired format | Organic synthesis, drug discovery, and materials design |

Table 2: Different types of prompt construction.

prompt to improve treatment recommendations. Similarly, CoI (Li et al., 2024b) organizes related research papers in a sequential chain to guide research idea generation, while ResearchAgent (Baek et al., 2024) incorporates background knowledge to create an iterative research agent that refines ideas.

**Iterative Feedback and Refinement** enables continuous improvement and adaptation of generated plans by using prompts to facilitate feedback and refinement. ResearchAgent (Baek et al., 2024) uses prompts to guide review agents in providing feedback to refine research ideas, while LogicSolver (Yang et al., 2022) prompts the LLM to not only solve mathematical problems but also explain the logical reasoning, enhancing transparency and interpretability of the planning process.

**Task Structuring and Multi-Agent Coordination** supports more complex planning by structuring tasks and coordinating actions across multiple agents or tools. Coscientist (Boiko et al., 2023) provides specific instructions for autonomous chemical experimental design and execution. It leverages four system prompts as commands that define the action space - 'GOOGLE', 'PYTHON', 'DOCUMENTATION', and 'EXPERIMENT'. ASA (Liu et al., 2024a) embeds the entire research cycle in the prompt for automated simulation, and VirSci (Su et al., 2024) organizes agents with role and task descriptions to collaboratively generate and refine research ideas. ChemCrow (Bran et al., 2024) tailors the prompt for managing chemistry-specific tools for tasks like materials design.

The above examples illustrate that prompt-based planning leverages carefully engineered prompts to trigger in-context learning, enabling scientific agents to generate logical, structured, step-by-step plans without extra fine-tuning. This approach embeds detailed context to guide decision-making in complex scientific tasks.

### 2.1.2 SFT-Based Planner

While prompt engineering and in-context learning offer zero-shot or few-shot planning capabilities for scientific agents, Supervised Fine-Tuning (SFT)-based planners enhance these capabilities by adapting a pre-training mechanism to specific scientific domains, as illustrated in Figure 3(b). The planning capability of SFT-based planners emerges from fine-tuning on **domain-specific planning trajectories**, which are curated datasets consisting of labeled input-output pairs. These pairs capture the step-by-step reasoning required for complex scientific tasks. For example, given a pre-trained planner with parameters $\theta$, SFT optimizes these parameters by training on domain-specific pairs $(x, y)$ derived from planning trajectories $D = (x_i, y_i)_{i=1}^{N}$. The objective function minimizes the negative log-likelihood:

$$\mathcal{L}_{SFT}(\theta) = -\mathbb{E}_{(x,y)} \sum_{t=1}^{T} log \ P_\theta(y_t|x, y_{<t}) \quad (1)$$

where $T$ denotes the planning step horizon and $y_{<t}$ represents previously generated planning steps. By training on such domain-specific data that pair complex scientific tasks with expert-level step-wise solutions, SFT-based scientific agent planners effectively bridge the "reasoning gap" observed in prompt-based methods, particularly for multi-step tasks like experimental design and hypothesis refinement.

For example, in drug discovery, DrugAssist (Ye et al., 2023a) utilizes an instruction-based dataset to fine-tune a planner for interactive molecule optimization. This training process enables the planner to internalize expert feedback and integrate it into the planning process, effectively creating a drug discovery agent. Similarly, ToolLLM (Qin et al., 2023) fine-tunes its planner on a specialized ToolBench dataset, improving its ability to invoke and interact with external APIs. By training on sequences of successful tool usage and interactions, the planner learns to generate plans that leverage external tools effectively, creating a tool-augmented scientific agent.

In summary, SFT-based planners provide a robust mechanism for aligning capabilities with the nuanced demands of scientific research. These planners are trained to replicate expert strategies by learning from annotated data that provides step-by-step solutions to complex scientific problems, enabling the planner to adapt and apply these strategies to new tasks. Examples like BioGPT (Luo et al., 2022), GatorTronGPT (Peng et al., 2023), and others show how fine-tuned planners can tackle diverse scientific tasks by internalizing the reasoning processes from expert-curated data.

### 2.1.3 RL-Based Planner

Reinforcement Learning (RL) plays a critical role in developing the planning capabilities of scientific agents by enabling them to autonomously refine decision-making strategies within complex scientific tasks (Rafailov et al., 2024). Unlike traditional planning systems, RL-based planners rely on feedback loops where agents learn to improve their actions based on rewards and penalties, for example, calculated from preference data as illustrated in Figure 3(c) by some reward functions. These planners are designed to receive positive reinforcement for desirable outcomes, such as accurate hypotheses or optimal experimental designs, and negative feedback for undesirable ones, such as logical errors. This learning process equips the planner with the ability to adapt over time, transcending the limitations of approaches like SFT.

The process of enhancing planning through RL is grounded in the agent's objective of maximizing cumulative rewards. Formally, this can be expressed as the optimization of a policy $\pi_\theta$ with respect to the reward function $r(y, x)$, which measures the quality of an action (or plan) $y$ generated from a given input $x$. Additionally, the agent's policy is regularized by a term involving the Kullback-Leibler (KL) divergence, ensuring that the updated policy does not deviate significantly from the original SFT-initialized policy. The formulation is:

$$J(\theta) = \mathbb{E}_{(x,y)\sim\pi_\theta}[r(y, x)] - \lambda \cdot K_l(\pi_\theta || \pi_{SFT}) \quad (2)$$

where $\pi_\theta$ is the planner being optimized, $r(y, x)$ quantifies the quality of the generated sequence $y$ given input $x$, and the KL-divergence $K_l$ ensures the updated policy remains close to the SFT-initialized behavior. Here, the planner learns adaptive strategies through an iterative trial-and-error process, allowing it to make more informed decisions across multiple steps, which is particularly
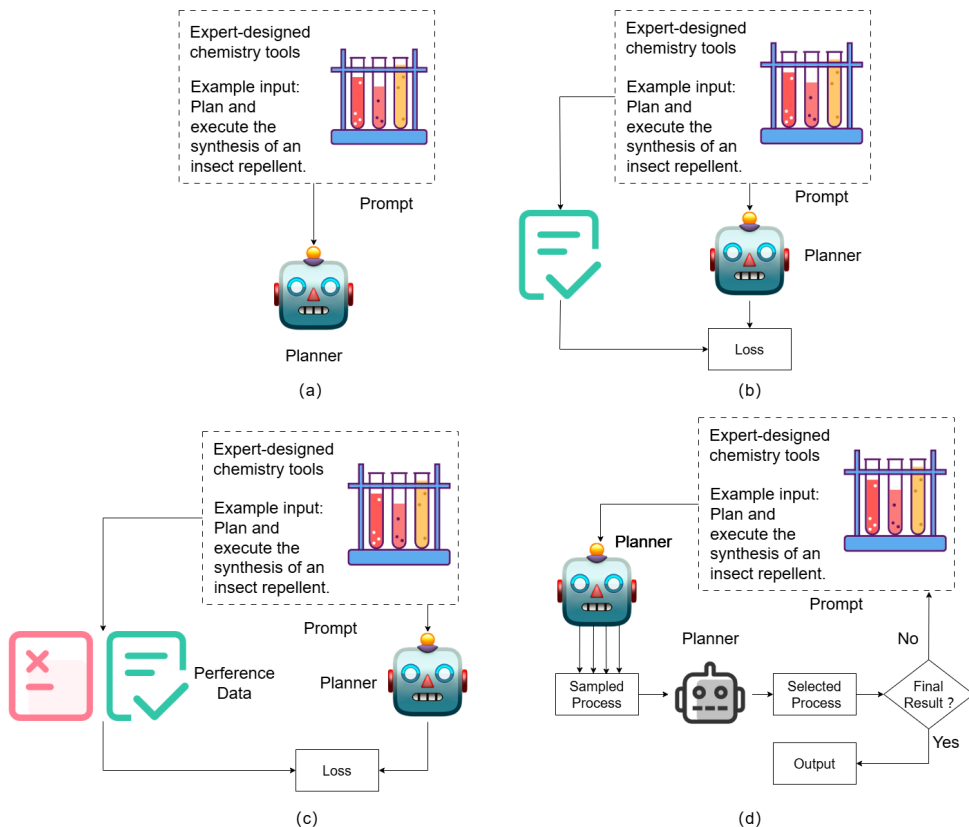
Figure 3: The types of planner in LLM-based scientific agents. (a) Prompt based planner; (b) SFT-based planner; (c) RL-based planner; (d) Process supervision based planner.

beneficial for complex scientific tasks requiring iterative refinement. We divide the studies of RL-based planners into three categories according to how they are designed.

**Iterative Refinement of Reasoning Paths.** To enhance the agent's reasoning and decision-making capabilities, RL-based planning incorporates iterative processes that refine the agent's reasoning paths. This is especially valuable in domains like mathematical problem-solving. For example, ReFT (Luong et al., 2024) uses Proximal Policy Optimization (PPO) combined with Chain-of-Thought (CoT) annotations, where an abundance of reasoning paths are automatically sampled given the question, and the rewards are naturally derived from the ground-truth answers. Similarly, CoT-Influx (Huang et al., 2024d) utilizes RL to optimize a coarse-to-fine pruning strategy that selects the most effective CoT examples for mathematical proofs. A multi-goal reward function is designed to measure the LLM loss, few-shot math reasoning effectiveness, and token length constraints.

Another important contribution comes from methodologies like STEP-DPO (Lai et al., 2024) and Flow-DPO (Deng and Mineiro, 2024), which apply Direct Preference Optimization (DPO) to the RL process. These techniques allow for the refinement of reasoning at the level of individual steps, treating each step as an element that can be optimized for better performance. This granular approach to optimization is crucial for improving the precision of reasoning in agents, making them more effective in complex scientific reasoning tasks.

**Optimizing Scientific Workflows and Design Spaces.** In scientific simulations and engineering design, RL-based planners are instrumental in optimizing intricate workflows and navigating large design spaces. The RL process in these contexts often involves multiple agents working together to optimize complex systems. For example, SciMARL (Scientific Multi-Agent Reinforcement Learning) (Bae and Koumoutsakos, 2022) demonstrates how a multi-agent RL framework for the discovery of wall models in large-eddy simulations, can identify optimal strategies for agents that perform actions, contingent on their information about the environment, and measures their performance via scenario-related scalar reward functions. In molecular discovery, RL can be employed to iteratively refine the design of de novo drug candidates. MolRL-

6

MGPT (Hu et al., 2024b) represents the problem of designing novel drug candidates as a cooperative Markov game consisting of multiple generative model agents with the scoring function as rewards. Similarly, Lutz et al. (2023) apply an RL-based approach combined with Monte Carlo tree search for protein design, creating complex protein nanomaterials with desired properties. These examples illustrate how RL can guide agents to explore and optimize design spaces, improving scientific outcomes across diverse fields.

**Integrating Human Feedback and Dynamic Environments.** For applications that require direct interaction with human expertise or adaptation to rapidly changing conditions, RL-based planners are specifically designed to integrate human feedback and adjust their strategies dynamically. For example, Barata et al. (2023) demonstrates how human preferences could be incorporated into a diagnostic AI for skin cancer by adjusting rewards and penalties based on expert-generated tables, thereby tailoring the system's performance to real-world clinical insights. For dynamic environments, Sauter et al. (2023) develops a meta-reinforcement learning algorithm for causal discovery. Their approach enables agents to construct explicit causal graphs with a reward based on the Structural Hamming Distance between the generated directed acyclic graph and the true causal graph, effectively guiding the agent toward more accurate causal models through optimal interventions. Additionally, Reinschmidt et al. (2024) applies RL to manage a magneto-optical trap in a cold atom experiment, where the system optimizes atom cooling by using a reward function that accounted for both the number of atoms and their average kinetic energy. This experiment underscores the robustness of RL-based planners in dealing with external disturbances.

These studies show that RL-based planners boost scientific agents' planning with well-designed rewards to refine decision-making and explore diverse solution paths for tasks like problem-solving, simulations, complex designing, and tasks incorporating human preference or in dynamic environments. They enable agents to autonomously enhance planning and reasoning, achieving precise and adaptable scientific intelligence over time and continuously improving overall performance.

### 2.1.4 Process Supervision Based Planner

Process supervision involves providing step-by-step feedback to Large Language Models (LLMs) during their reasoning or generation process, rather than only evaluating the final outcome. In recent studies, this technique has been employed to enhance the planning and reasoning capabilities of scientific LLMs. For instance, systems like Marco-o1 (Zhao et al., 2024a) integrate Chain-of-Thought (CoT) fine-tuning with Monte Carlo Tree Search (MCTS) and reflective mechanisms to explore multiple reasoning paths, while SCoRe (Self-Correction with Reinforcement Learning) (Kumar et al., 2024) leverages multi-turn online reinforcement learning on self-generated correction traces to continuously refine reasoning. Additional improvements include methods like V-STaR (Hosseini et al., 2024), which trains a verifier using Direct Preference Optimization (DPO) to select the best candidate among outputs, and OmegaPRM (Luo et al., 2024), which employs a divide-and-conquer strategy with MCTS to gather process supervision data that train Process Reward Models and enhance mathematical reasoning. These innovations not only strengthen the core reasoning process of LLMs, but also establish the foundation for process supervision based planners in LLM-based scientific agents, where similar feedback mechanisms are used to iteratively refine and optimize complex scientific hypotheses, as illustrated in Figure 3(d).

Building on these advances in process supervision for scientific LLMs, similar principles are adapted to the design of LLM-based scientific agents, yielding planning architectures that dynamically integrate automated hypothesis generation with domain-specific evaluative feedback. For example, ChemReasoner (Sprueill et al., 2024) leverages an LLM-driven planner to systematically navigate the expansive chemical space. In this framework, the LLM constructs a hierarchical search tree where each node embodies a distinct hypothesis generated through "query plans" that include catalyst type, inclusion/exclusion criteria, and relational operators. The planner then uses quantum-chemical feedback—derived from atomistic simulations evaluating adsorption energies, reaction energy barriers, and structural stability—to assign rewards that prune unpromising pathways and iteratively refine the hypothesis space. This dual-loop mechanism effectively guides the exploration toward energetically favorable catalysts.

Similarly, the Scientific Generative Agent (SGA) (Ma et al., 2024a) employs a bi-level optimization framework to enhance planning capabilities for scientific discovery. At the outer level, the LLM

functions as a strategic planner, generating discrete hypotheses and experimental designs while dynamically adjusting its query prompts based on past simulation results. In parallel, the inner level leverages differentiable simulations to optimize continuous parameters—such as physical constants or molecular coordinates—providing gradient-based feedback that informs subsequent hypothesis refinement. By balancing exploitation (refining known promising designs) and exploration (venturing into novel solution spaces) through controlled temperature tuning, SGA's planning cycle adapts to emergent data and uncertainties, thereby increasing both robustness and innovation in scientific outcomes.

Together, these strategies illustrate that embedding detailed, domain-specific feedback into the planning process empowers scientific agents to engage in continuous hypothesis refinement, adaptive experiment design, and iterative plan optimization. Such dynamic planning capabilities significantly enhance the agents' efficiency, accuracy, and adaptability in addressing complex scientific challenges.

### 2.1.5 Discussion

In summary, the Planner component—comprising prompt-based, SFT-based, RL-based, and process supervision-based approaches—serves as the central controller in LLM-based scientific agents. These planners are crucial for enabling scientific agents to translate high-level scientific queries into actionable plans by decomposing tasks, integrating domain-specific knowledge, and coordinating interactions with specialized tools. Scientific agents rely on logical, structured planning to ensure that experimental protocols and hypothesis testing are carried out methodically. However, despite advances in these approaches, challenges remain as shown in Table 1. For example, prompt-based planners are highly sensitive to prompt quality, leading to inconsistent scientific outputs; SFT-based planners require extensive, high-quality domain-specific datasets that are often costly to curate; RL-based planners struggle with designing robust reward functions and managing computational costs critical for scientific exploration; and process supervision-based planners, while promising for iterative refinement, demand complex and as-yet non-standardized feedback mechanisms.

Looking ahead, there are several promising directions for future research in the context of scientific agents. First, designing efficient surrogate models and robust reward mechanisms could reduce the computational burden associated with RL-based planning for scientific agents, making them more practical for real-world scientific problems. Second, integrating automated prompt optimization and self-supervised feedback could enhance the reliability and scalability of prompt-based and process supervision-based planners within scientific agents, leading to more consistent and accurate scientific outputs. Finally, establishing standardized evaluation benchmarks and cross-domain interface protocols will be essential for tracking progress and ensuring that future LLM-based scientific agents are both effective and ethically sound. These efforts will collectively contribute to building more autonomous, transparent, and efficient scientific agents capable of driving rapid, reproducible, and innovative scientific discovery.

### 2.1.6 Single-agent vs. Multi-agent Planner

As we note under Figure 1, the planner could be implemented in a single-agent fashion — where one LLM handles all planning, reflection, and verification functions — or in a multi-agent manner, with specialized agents distributed to execute these distinct tasks.

Single-agent planners integrate all core functions—task planning, reflection, memory access, and tool use—into a single unified module. This monolithic design simplifies system architecture and debugging, making it well-suited for applications where the scope of tasks is limited, or tight integration between components is essential. In early scientific agent systems, such as Coscientist (Boiko et al., 2023) and ChemCrow (Bran et al., 2024), a single LLM-based planner would orchestrate all operations, providing a streamlined approach that is easier to implement and manage. However, this simplicity can become a bottleneck when addressing complex scientific problems that require specialized subtasks to be executed concurrently or with varying degrees of autonomy.

In contrast, multi-agent planners, exemplified by recent developments like Google's AI co-scientist (Gottweis et al., 2025), distribute these responsibilities among specialized agents (see Appendix A for detailed introduction). In such architectures, distinct agents may be assigned to generate hypotheses, perform critical reflection, rank and refine ideas, or even manage meta-reviews. This division of labor enables a "generate, debate, and evolve" framework where each agent focuses on a specific function, enhancing overall system flex-
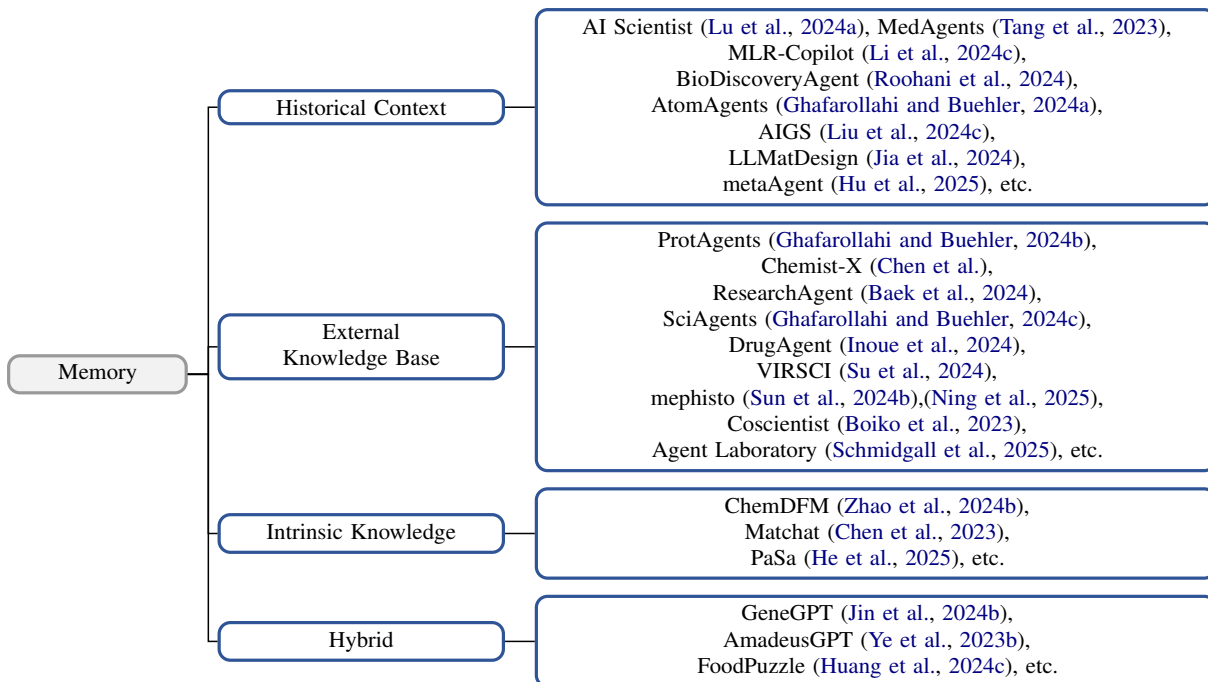
Figure 4: Taxonomy of the memory mechanism of science agents.

ibility and scalability. Empirical results from the AI co-scientist framework indicate that this modular approach can significantly accelerate discovery processes—for example, by reducing hypothesis generation timelines from weeks to days—and improve the novelty and accuracy of research outputs. On the one hand, the increased complexity of multi-agent systems demands robust communication protocols and coordination strategies to manage potential inter-agent conflicts and ensure coherent output. On the other hand, they offer enhanced performance in tackling multifaceted, interdisciplinary challenges typical in scientific research.

In summary, the choice between single-agent and multi-agent planners depends largely on the task complexity, the need for specialization, and the desired scalability of the system. For routine or narrowly defined problems, single-agent planners may suffice, whereas multi-agent planners are better suited for advanced scientific discovery where dynamic, specialized collaboration is key.

## 2.2 Memory

Memory in LLM-based scientific agents extends beyond simple context retention, enabling long-term accumulation of research findings, iterative hypothesis refinement, and cross-project continuity. By mirroring the cognitive processes of human scientists, these agents maintain detailed historical context, integrate domain-specific external knowl-

edge, and leverage intrinsic model capabilities to ensure that each experiment or literature insight informs future decisions. We categorize these memory mechanisms into Historical Context, External Knowledge Base, and Intrinsic Knowledge—three facets that collectively address the timeline-driven nature of scientific inquiry, the breadth of specialized data sources, and the deep, model-level understanding required for advanced tasks. While not mutually exclusive, each category highlights a distinct dimension of how scientific agents store and utilize information to reproduce results, accumulate evidence, and push the boundaries of autonomous research. We compare the three mechanisms in Table 3, and list the related studies in Figure 4.

### 2.2.1 Historical Context

Historical context—often termed conversational or short-term memory—is vital for scientific agents to maintain continuity and iterative progress in research workflows. Unlike general agents that merely hold transient dialogue, scientific agents accumulate and leverage past interactions, experimental outcomes, and reasoning steps to refine hypotheses and improve experiment designs over time. This robust memory enables them to mimic the cumulative nature of scientific inquiry, ensuring each cycle of analysis builds on previous insights and supports reproducible results. Figure 5 illustrates how historical context underpins the iterative

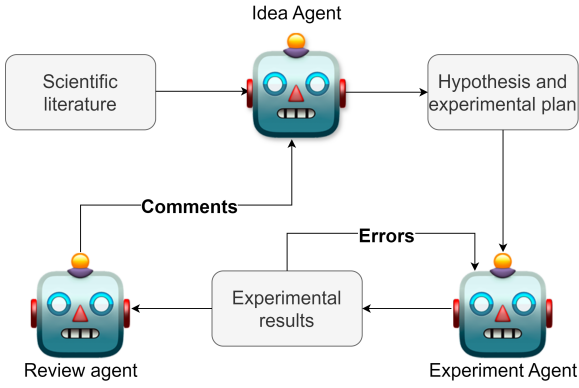| Type | Methodology | Strengths and Limitations | Typical Use Cases |
|------|-------------|---------------------------|-------------------|
| Historical Context | Maintains conversational logs or iterative action sequences; stores previous interactions and experimental outcomes | ✓ Enables a coherent and iterative refinement<br>✓ Supports dynamic adaptation<br>✗ Limited by the model's context window<br>✗ Implicit logging can make explicit retrieval challenging | Iterative hypothesis refinement, tracking multi-turn research sessions, and adapting strategies based on previous interactions. |
| External Knowledge Base | Accesses curated external sources such as literature databases and structured knowledge graphs | ✓ Expansive, up-to-date domain-specific information<br>✓ Reasoning in external, validated research<br>✗ Integration can be complex<br>✗ Dependent on the quality and update frequency of external sources | Comprehensive literature reviews, retrieving domain-specific data. |
| Intrinsic Knowledge | Represents the inherent capabilities and information embedded in the LLM from pre-training and fine-tuning. | ✓ Provides a robust foundation for general language understanding and scientific reasoning<br>✓ Immediately available<br>✗ May become outdated over time<br>✗ Limited by the scope and recency of the training data | General scientific reasoning, initial hypothesis generation, and immediate general foundational tasks. |

Table 3: Comparison of different memory types.



Figure 5: A simple process of scientific agents using historical context (e.g.,comments provided by the Review Agent, and errors in each round of experiments).

refinement process in scientific agents.

Several frameworks highlight the importance of historical context, albeit with varying implementations. For instance, AI Scientist (Lu et al., 2024a) utilizes historical context by iteratively developing ideas and adding them to a growing archive, mimicking the cumulative knowledge building in the scientific community. Similarly, MedAgents (Tang et al., 2023) emphasizes iterative discussions to reach consensus, where the progression of arguments within the dialogue itself forms the historical context. Similarly, MLR-Copilot (Li et al., 2024c) and BioDiscoveryAgent (Roohani et al., 2024) utilize feedback from previous rounds' experiments to refine their subsequent steps, directly incorporating past results into the ongoing process. AtomAgents (Ghafarollahi and Buehler, 2024a) provides a more structured approach by defining dedicated "core

memory" and "tool memory" modules to store conversations between agents and tool interactions, ensuring readily accessible historical data throughout problem-solving. AIGS (Liu et al., 2024c) further illustrates this with its "Pre-Falsification" phase that relies on multi-turn logs of iterative exchanges between agents, explicitly using these logs as the history context to refine proposals. Recent works further highlight sophisticated uses of history context. LLMatDesign (Jia et al., 2024) incorporates "self-reflection" on previous decisions, allowing the agent to rapidly adapt to new tasks and conditions in a zero-shot manner. The metaAgent (Hu et al., 2025), an embodied intelligent agent for electromagnetic space, leverages a "multi-agent discussion mechanism" in its "cerebrum" part. This mechanism inherently relies on the conversations and interactions between specialized agents.

Despite varied implementations, historical context is essential for maintaining a continuous record of interactions. It enables agents to iteratively refine their approaches based on past successes, failures, and external inputs. Whether through explicit memory modules or implicit conversational flow, this iterative process guides future actions, ensuring dynamic adaptation and coherent execution in the scientific discovery process.

### 2.2.2 External Knowledge Base: Augmenting Agent Capabilities with Broad Scientific Knowledge

External knowledge bases (KBs) are essential for scientific agents, providing a curated repository of up-to-date, domain-specific information that ex-

tends beyond the static training data of LLMs. These KBs are not merely supplemental—they are deeply integrated into the agent's reasoning process, enabling it to retrieve, synthesize, and connect complex scientific concepts. This external integration is critical for tasks that demand in-depth domain expertise and comprehensive literature awareness. By systematically incorporating external knowledge, scientific agents can enhance hypothesis generation, experimental design, and data analysis, ensuring that their outputs remain current, robust, and contextually relevant. Figure 6 illustrates this process.

A prominent approach remains leveraging **scientific literature** as an external KB. ProtAgents (Ghafarollahi and Buehler, 2024b) and Chemist-X (Chen et al.) both employ Retrieval-Augmented Generation (Lewis et al., 2020) with literature databases, allowing agents to ground their reasoning in existing research. ResearchAgent (Baek et al., 2024) takes a more structured approach, building an "entity-centric knowledge store" from literature co-occurrences to capture underlying relationships and facilitate cross-pollination of ideas. Agent Laboratory (Schmidgall et al., 2025) illustrates utilization of literature through the arXiv API, enabling agents to retrieve, summarize, and generate papers.

Beyond literature, **knowledge graphs (KGs)** emerge as another significant type of external KB. SciAgents (Ghafarollahi and Buehler, 2024c) explicitly uses large-scale ontological KGs to organize scientific concepts, ensuring generated scientific hypotheses are rooted in interconnected scientific concepts. DrugAgent (Inoue et al., 2024) uses a Knowledge Graph Agent to extract drug-target interaction information, demonstrating the use of targeted KGs for specific domains. By adopting the KnowledgeBank module from AgentScope (Gao et al., 2024b), VirSci (Su et al., 2024) makes scientist agents' profiles embed into the author knowledge bank, through which agents can quickly access and familiarize themselves with other initialized agents' information.

Expanding beyond traditional literature and KGs, Coscientist (Boiko et al., 2023) demonstrates the power of integrating diverse resources as external KBs. It mainly intergrates Web Searcher module and Documentation search module to enable the agent browse the internet and relevant documentation for experiments in next period. In geospatial domain, an autonomous geospatial data re-
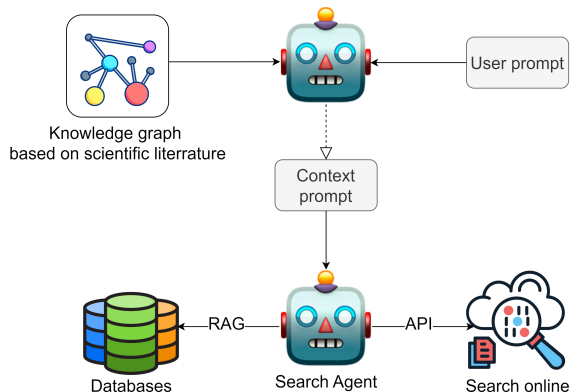


Figure 6: A simple process of scientific agents using external knowledge base.

trieval framework (Ning et al., 2025) manages a pre-defined and scalable list of data sources like OpenStreetMap and US Census data, highlighting a more curated external data approach for specific tasks. In Astronomical domain, mephisto (Sun et al., 2024b) adopts a dynamically updated external knowledge base through a learning system where domain knowledge is extracted and then validated, suggesting a form of continuously learning and evolving external knowledge.

These examples demonstrate a clear trend: scientific agents are significantly enhanced by access to diverse external KBs. The KB type varies greatly based on the task, ranging from broad scientific literature and web resources to specialized KGs, curated datasets, etc. The functional overlap across the approaches lies in their ability to provide agents with access to a wider and more current scope of information than their intrinsic knowledge. Complementarily, external KBs ground agent reasoning in established scientific knowledge, promote validity and novelty by referencing existing works and data, and enable interaction with real-world tools and information sources. Implementation-wise, we see a spectrum from RAG-based retrieval from unstructured text to direct querying of structured KGs and databases, API integrations, and web browsing, reflecting increasingly sophisticated strategies for knowledge integration and utilization to empower scientific agents.

### 2.2.3 Intrinsic Knowledge: Leveraging Pre-trained and Fine-tuned LMs

In the context of scientific agents powered by Large Language Models (LLMs), intrinsic knowledge of LLMs serves as the foundational cognitive bedrock. This refers to the inherent capabilities and infor-

mation that the LLM itself embodies, meticulously cultivated during its pre-training phase on massive and diverse corpora, crucially including scientific literature, datasets, and domain-specific knowledge. This intrinsic knowledge is further refined through task-specific fine-tuning. For a scientific agent, this isn't merely passive data storage; it's the very source of an agent's reasoning faculties, natural language competency, and fundamentally, its foundational scientific understanding. The intrinsic knowledge, therefore, empowers a scientific agent to operate effectively within scientific contexts, providing the essential base for scientific reasoning, comprehension of scientific language, and the broad scientific literacy required to function as an autonomous scientific explorer and problem-solver.

Several studies emphasize the importance of enhancing the intrinsic knowledge of LLMs for scientific agents through specialized training. ChemDFM (Zhao et al., 2024b) pioneers this approach by developing a domain-specific LLM pre-trained on a massive chemical literature and textbook corpus and further fine-tuned with chemical instructions. This directly injects chemical expertise into the model's core knowledge. Matchat (Chen et al., 2023) takes a fine-tuning route, enhancing Llama2-7B with structured material knowledge data, demonstrating the efficacy of incorporating domain-specific structured information to improve model performance in materials science. PaSa (He et al., 2025) focuses on academic paper search, utilizing reinforcement learning with a synthetic dataset of academic queries and papers to optimize an LLM agent for search task.

Building upon these strategies, recent works further explore diverse avenues for enriching intrinsic knowledge. ProLLaMA (Lv et al., 2024) introduces efficiency into the fine-tuning process for Protein Language Models by employing Low-Rank Adaptation (Hu et al., 2022). This method improves the efficiency of protein learning during fine-tuning, demonstrating advancements in making specialized model training more resource-effective. Moreover, Tx-LLM (Chaves et al., 2024) presents a generalist large language model for therapeutics, fine-tuned from PaLM-2 (Anil et al., 2023). Tx-LLM distinguishes itself by being trained on an extensive collection of 709 datasets, encompassing 66 tasks across the drug discovery pipeline. In contrast to single domain fine-tuning, NatureLM (Xia et al., 2025) adopts a multi-domain pre-training approach. Pre-trained on data from multiple scientific domains that include small molecules, materials, proteins, DNA and RNA, NatureLM aims to offer a unified and versatile model applicable across various scientific applications.

These examples highlight a critical strategy : tailoring the LLM's intrinsic knowledge to the specified scientific domain. The functional overlap is clear – all approaches aim to improve the LLM's base capabilities for scientific reasoning and task execution within specific fields, whether it be chemistry, materials science, academic search, protein science, therapeutics, or broadly across multiple scientific disciplines. Complementarily, intrinsic knowledge provides the bedrock for the agent's intelligence, enabling it to effectively process historical context and utilize external knowledge. Implementation approaches differ significantly, ranging from full domain-specific pre-training (ChemDFM, NatureLM) to targeted fine-tuning with structured data (Matchat, Tx-LLM) or reinforcement learning (PaSa), and including techniques for efficient fine-tuning (ProLLaMA). These diverse techniques underscore the importance of carefully shaping the LLM's intrinsic knowledge and demonstrate the expanding LLMs available for researchers to create scientifically intelligent agents.

### 2.2.4 Discussion

Memory in LLM-based scientific agents is implemented via three interrelated mechanisms: history context, external knowledge bases, and intrinsic knowledge. History context enables agents to maintain conversational coherence and iterative refinement by retaining and recalling prior interactions, emulating the cumulative nature of human research. External knowledge bases expand the agent's informational scope by integrating up-to-date and domain-specific data, allowing for the retrieval, synthesis, and contextualization of complex scientific concepts. Meanwhile, intrinsic knowledge enables agents to apply core scientific reasoning from the outset, serving as the bedrock for advanced, context-rich memory layers.

Despite their complementary roles, current memory mechanisms face several limitations. Many approaches—especially those using textual memory-suffer from scalability issues and information loss since context windows are limited. Parametric methods, while more efficient, often lack interpretability and require extensive fine-tuning. Moreover, external knowledge integration remains brittle in dynamically changing domains, leading to po-

tential mismatches or outdated retrievals. Recent studies (Xu et al., 2025; Zeng et al., 2024) emphasize the need for more adaptive, self-organizing memory systems that can dynamically link and update stored information.

Future research should focus on developing hybrid memory models that combine the benefits of both parametric and textual representations, or exploring the synergistic relationships between different memory types and investigate novel hybrid approaches to optimize their collective performance in automated scientific discovery (e.g., GeneGPT (Jin et al., 2024b), AmadeusGPT (Ye et al., 2023b), FoodPuzzle (Huang et al., 2024c)). Further, integrating robust metadata learning and external knowledge graphs—as explored in recent works like Hatalis et al. (2023)—could enhance retrieval accuracy and contextual grounding. Additionally, improved lifelong learning techniques and efficient forgetting mechanisms are essential to mitigate memory overload and maintain performance over extended research cycles.

## 2.3 Tool Set

In this section, we introduce the tool sets employed by scientific agents, as illustrated in Figure 7. While LLMs demonstrate robust problem-solving capabilities for general tasks and foundational scientific inquiries, they often encounter limitations when addressing advanced scientific challenges, particularly those in STEM-related domains, due to insufficient domain-specific expertise and computational resources. The tool set extends the LLM's capabilities beyond natural language processing by enabling real-time data retrieval, precise code execution, domain-specific scientific computation, and rigorous experimental simulation. This tight integration allows scientific agents to access accurate, up-to-date information, perform computationally intensive analyses, and process data in specialized modalities—capabilities that are essential for simulating and validating experiments. Consequently, these tool sets serve not just as supplementary resources but as a core component of the agent's architecture, fundamentally enhancing its scientific reasoning, reliability, and adaptability in complex research environments.

Based on the functional types of tools, we categorize existing scientific agents tool sets into two categories: (1) Tool sets based on APIs and code libraries, and (2) Tool sets based on simulators or emulation platforms. The subsequent section will present recent advancements in each category.

### 2.3.1 Tool sets based on APIs and code libraries

Tool sets based on APIs and code libraries aim to extend the knowledge boundaries and computational capacities of LLMs in scientific tasks. These tool sets encapsulate domain-specific knowledge bases and specialized algorithm libraries into standardized functional interfaces, thus enabling LLMs to transcend the limitations imposed by the timeliness and domain depth of their training data, as well as computational limitations inherent in LLMs. This category encompasses both pre-existing general-purpose tools, discipline-specific scientific tools, and novel tools synthesized by researchers using generative methods.

Simpler tool sets encompass basic components such as search engines or database query modules. For instance, MAPI-LLM (Jablonka et al., 2023) leverages LLMs to retrieve information from the Materials Project API (MAPI) Reaction-Network package and Google Search, addressing user queries about chemical materials. The LLM employs Chain-of-Thought prompting to translate natural language queries into structured API calls, allowing users to retrieve material properties and execute complex searches through conversational interfaces. Similarly, ClimateGPT (Thulke et al., 2024) integrates a retrieval mechanism to access a curated collection of climatological research reports and peer-reviewed papers, thereby enhancing response accuracy in climate science applications.

While the aforementioned examples focus on elementary query tools, the following case demonstrates the integration of sophisticated and multifunctional APIs, which significantly augment the capabilities of scientific agents in handling diverse scientific tasks. These advanced tool sets not only grant access to domain-specific data repositories but also enable intricate computational workflows and analytical operations, thereby pushing the boundaries of LLMs in scientific reasoning.

Mathematics is frequently regarded as the foundational discipline underpinning numerous scientific domains. Tora (Gou et al., 2024) integrates Python libraries such as SymPy, SciPy, and CVXPY into natural language reasoning frameworks, demonstrating significant performance improvements for open-source LLMs across multiple mathematical reasoning benchmarks. In the disciplines of chemistry and materials science, Chem-
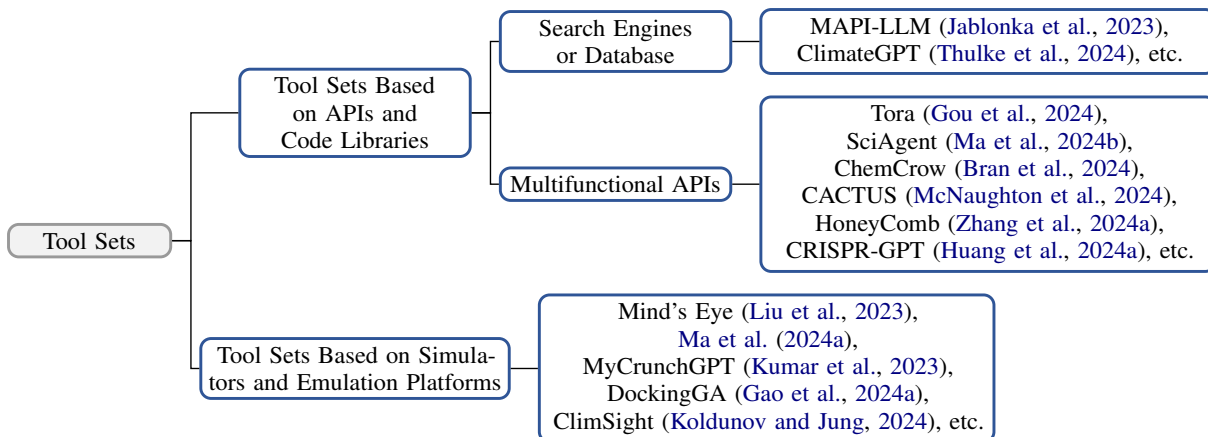
Figure 7: Taxonomy of the tool sets of scientific agents.

Crow (Bran et al., 2024) deploys an extensive tool set comprising 18 expert-designed tools that support functions such as molecular property queries, reaction prediction, and experimental synthesis planning. The integration of these tools empowers LLMs to autonomously design and execute complex chemical workflows in organic synthesis, drug discovery, and materials design. Similarly, CACTUS (McNaughton et al., 2024) enhances cheminformatics capabilities by integrating tools from open-source packages such as RDKit, enabling quantitative estimation of drug-likeness and pharmacokinetic properties for chemical compounds provided in SMILES. The HoneyComb (Zhang et al., 2024a) framework constitutes a comprehensive materials science framework, integrating MatSciKB (Materials Science Knowledge Base) with a ToolHub. MatSciKB gathers structured knowledge from peer-reviewed literature, while the ToolHub incorporates search engines, Python interpreters, and domain-specific APIs constructed via an inductive tool construction methodology. In the domain of biology, CRISPR-GPT (Huang et al., 2024a) synergizes Google Search, Primer3, Broad Institute's gold-standard guideRNA library, the CRISPRPick tool set, and scholarly databases. This integration enables researchers to select the most suitable CRISPR systems and to design experimental protocols for genome-editing workflows. SciAgent (Ma et al., 2024b) introduces a methodology to generalize LLMs' mathematical tool utilization to other scientific domains. Researchers initially generated a mathematical tool set via a cross-retrieval strategy, subsequently developing a human-validated and refined multi-domain tool set based on the SciToolBench dataset. This comprehensive tool set encompasses five disciplines: Mathematics, Physics, Chemistry, Finance, Electrical Engineering, and Computer Science.

The integration of these tool sets into scientific agent systems has enabled researchers to enhance LLMs' capacities in experimental planning and numerical prediction within many scientific domains. This modular, extensible integration strategy has proven effective in mitigating LLMs' inherent limitations in domain expertise and computational precision. Nevertheless, challenges such as non-standardized interfaces, limited tool diversity, and the complexity of tool generation are currently hindering the broader adoption of such tool sets.

### 2.3.2 Tool sets based on simulators and emulation platforms

Tool sets based on simulators and emulation platforms provide specialized, domain-specific tools for scientific agents, enabling them to simulate experimental procedures and validate results. By translating natural language instructions into executable simulation codes or parameterized control signals using LLMs, these tool sets facilitate deep integration with experiment workflows. Often, they are tightly coupled with the planning process, ensuring correct parameterization and validation throughout simulations or lab automation steps, and especially valuable in complex research tasks.

At present, this class of tool sets is employed most frequently in physics-related scientific agents. Mind's Eye (Liu et al., 2023) employs the MuJoCo physics engine to simulate real-world physical scenarios. The language model converts natural language text into rendering codes, with simulation results iteratively incorporated into subsequent inputs, thereby facilitating physics-based reasoning.

Similarly, Ma et al. (2024a) utilize physics simulators as experimental platforms on which LLMs generate scientific hypotheses and perform reasoning. The simulator provides observational feedback and enables differentiable optimization of continuous parameters, thereby achieving validated results in constitutive law discovery and molecular design tasks. MyCrunchGPT (Kumar et al., 2023) integrates a suite of software components, including DeepONet surrogates and the computational fluid dynamics (CFD) simulator Nektar++ to optimize 2D NACA airfoils in aerodynamic design. The LLM employs DeepONet for flow field computations during the process of airfoil optimization, and the validity of the results is confirmed through high-fidelity simulations. In the domain of chemistry, DockingGA (Gao et al., 2024a) utilizes molecular docking simulations to facilitate the generation of molecules that exhibit target-specific binding affinities. The docking scores between the generated molecules and biological targets function as reward signals, thereby enabling the refinement of molecular synthesis. In the context of climatology, ClimSight (Koldunov and Jung, 2024) integrates geospatial databases and the AWI Climate Model, a global climatology simulation framework, to assess the climate impacts of specific activities, such as agricultural practices.

The integration of simulation tool sets is a solution to the limitations of LLMs in understanding physical laws and reasoning about dynamic processes. This enhances computational accuracy and validity for complex problems. However, the practical adoption of such tool sets remains constrained by the high computational costs and temporal overheads of precision simulators. In addition, the proficient utilization of simulators and the accurate generation of their corresponding parameters also pose significant challenges for LLMs.

### 2.3.3 Discussion

Recent studies have shown that the incorporation of scientific tool sets into agent systems leads to substantial enhancements in LLMs' planning, reasoning, computational, and execution capabilities for scientific tasks. Tool sets based on APIs and code libraries address limitations in domain knowledge and computational power by encapsulating specialized algorithms and knowledge bases. This separation allows scientific agents to decouple high-level reasoning from raw numerical operations, enabling them to prioritize strategic planning and or-

chestrate complex tool usage. Simultaneously, tool sets based on simulator and emulation platforms integrate experimental simulations with natural language reasoning, augmenting the agents' ability to manage and solve intricate, multi-step scientific workflows with improved precision and reliability.

However, several limitations persist in current tool integration studies from the scientific perspective. Many systems still rely on pre-defined tool sets and static, well-documented repositories, which restrict scalability and adaptability in dynamic research environments. For example, benchmarks like ShortcutsBench (Shen et al., 2025) reveal that even state-of-the-art systems struggle with managing API dependencies and adapting to frequently updated external services—challenges that are particularly acute in rapidly evolving scientific domains such as computational biology and materials science. In addition, high subscription costs for some API services, along with persistent challenges in error handling, security, and reproducibility, continue to impede the deployment of robust LLM-based agents in rigorous scientific research.

Looking ahead, future research must develop autonomous, self-adaptive tool generation frameworks that leverage middleware layers to seamlessly integrate diverse functionalities at runtime. Promising strategies, as highlighted in recent works such as Shen et al. (2025); Gu et al. (2024), suggest that dynamic middleware-based solutions can adapt to real-time changes in scientific environments. Moreover, standardizing API design and documentation, enhancing automated error detection and recovery mechanisms, and creating comprehensive, dynamic benchmarks tailored for scientific applications will be pivotal.

## 3 General-purpose vs. Scientific Agent

The above section shows the module design for scientific agents. Different from scientific agents (e.g., AI Co-Scientist (Gottweis et al., 2025)) that specialize in research workflows, general-purpose agents (e.g., Manus (Manusai.ai, 2025)) are designed for broad adaptability across user-defined tasks. While they may share foundational LLM technology, their planning, memory strategies, tool integrations, and reasoning approaches differ significantly. This section outlines key technical distinctions that necessitate dedicated scientific agent design. Table 4 lists the key different features. Noting that current scientific agents are still in the early stage, perhaps

| Aspects | General-purpose Agents | Scientific Agents |
|---|---|---|
| Planning and Task Management | - Heuristic or reactive planning<br>- Flexible, goal-driven methods<br>- Not aligned with scientific methodology | - Logical, structured, hierarchical planning<br>- Long-horizon research projects<br>- Mirrors the scientific method |
| Memory and Knowledge Integration | - Ephemeral, context-limited storage<br>- Typically single-session or ad-hoc<br>- Minimal cross-project continuity | - Persistent, structured memory<br>- Accumulates data and insights across multiple experiments<br>- Enables reproducibility and long-term progression |
| Tool Utilization and Integration | - Plugin-based for a wide variety of tasks<br>- Minimal domain-specific parameterization | - Specialized, domain-specific tools<br>- Deep integration for experiment workflows |
| Domain-Specific Reasoning and Collaboration | - Mostly direct, goal-focused<br>- Single-agent or loosely multi-agent<br>- Often relies on user feedback to catch errors | - Iterative, hypothesis-driven logic<br>- Integrates domain rules and scientific practices<br>- Multi-agent debate and consensus-building<br>- Rigorous statistical checks, error bounds |

Table 4: Comparison between general-purpose agents and scientific agents.

no single agent system has yet achieved all the features in the table, but this is their trend, due to the high logic, structured content, long-term retention, professionalism, low tolerance for error, and reproducibility of the scientific field.

## 3.1 Planning and Task Management

**General-purpose agents** often use *heuristic* or *reactive* planning approaches (e.g., ReAct (Yao et al., 2023), plan-then-execute (Zhang et al., 2025)), adjusting actions based on intermediate results to maintain flexibility. Although multi-agent designs like Manus (Manusai.ai, 2025) allow broader task delegation, they generally lack **formal structure** that enforces scientific methodology over long, complex research phases.

**Scientific agents**, by contrast, implement *structured, hierarchical* planning aligned with the scientific method (Schmidgall et al., 2025; Gottweis et al., 2025). Systems such as BioPlanner (O'Donoghue et al., 2023) systematically translate scientific goals into reproducible protocols, and AI Co-Scientist (Gottweis et al., 2025) uses **parallel multi-agent planning** to handle literature review, hypothesis generation, and ranking in tandem. This **logical framework** ensures that each step—hypothesis, experiment, analysis—proceeds in a coherent, method-driven sequence.

## 3.2 Memory and Knowledge Integration

**General-purpose agents** typically rely on *ephemeral* memory, constrained by context windows or retrieval-augmented generation (RAG) (Park et al., 2023). Tools like AutoGPT (Yang et al., 2023) may store short-term notes (scratchpads), but they rarely support **long-term accumulation** of information. As a result, knowledge retention is ad-hoc, depending on frequent web queries instead of persistent internal structures.

**Scientific agents** emphasize **persistent memory** that evolves throughout extended research projects. AI Co-Scientist (Gottweis et al., 2025) maintains a *shared memory store* of intermediate results, accessible to specialized sub-agents for a coherent, team-like workflow. LLaMP (Chiang et al., 2024) and Agent Laboratory (Schmidgall et al., 2025) integrate structured domain databases, enabling **cumulative knowledge** across multiple projects. This large-scale retention preserves experimental histories, fosters reproducibility, and supports the long timelines inherent in scientific investigation (see also Lu et al. (2024a)).

## 3.3 Tool Utilization and Integration

**General-purpose agents** usually adopt a *plugin-based* model (e.g., Toolformer (Schick et al., 2023), HuggingGPT (Shen et al., 2024a)), calling APIs like web search or Python execution as needed. These integrations are generic, supporting various tasks but lacking specialized simulation or experiment workflows.

**Scientific agents**, on the other hand, require **deeply integrated tools** for simulations, experiment orchestration, and data analysis. For example, specialized modules in ProtAgents (Ghafarollahi and Buehler, 2024b) handle complex computational biology tasks, while chemistry-focused frameworks (e.g., ChemCrow (Bran et al., 2024)) support reaction prediction and laboratory automation. Crucially, these tools are not just *invoked*—they are part of the **scientific planning loop**, ensuring that parameters, methods, and validations conform to domain standards.
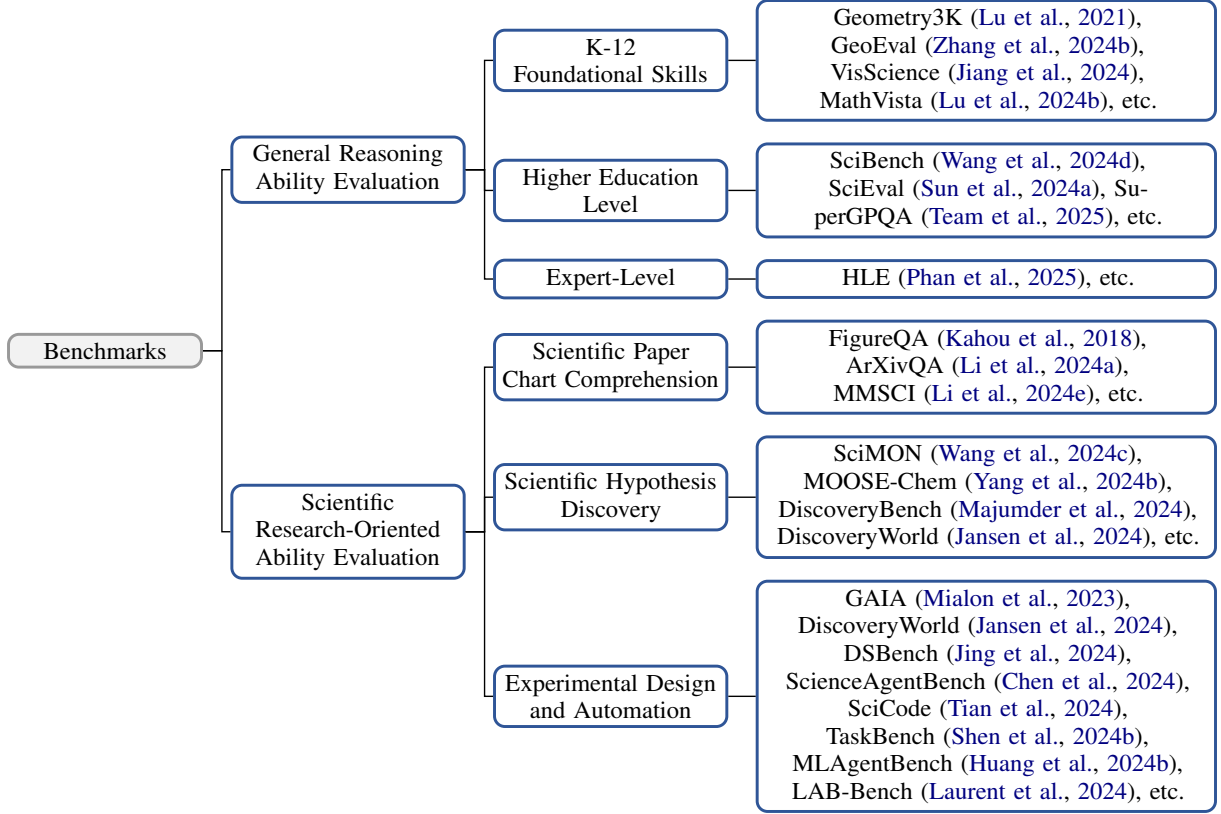
Figure 8: Taxonomy of the LLM-based science agents evaluation benchmarks.

## 3.4 Domain-Specific Reasoning and Collaboration

**General-purpose agents** focus on achieving single-user goals, often lacking *built-in verification*. Self-reflection may occur, but thorough validation (e.g., statistical checks, error bounds) is not typically included.

**Scientific agents** implement **validation and reproducibility** measures to ensure robust outputs. Multi-agent debate (Su et al., 2024) allows for hypothesis refinement via critical discussion, and AI Co-Scientist (Gottweis et al., 2025) employs parallel hypothesis evaluation, discarding flawed ideas early. By incorporating statistical analyses, error checking, and domain-specific constraints, scientific agents **uphold reliability and reproducibility**—core requirements of scientific research.

## 4 Benchmark

Benchmarks are basic solutions for evaluating the efficacy of LLM-based scientific agents, ensuring their capability to handle the multifaceted demands of scientific research. They are designed to measure various aspects of these agents' performance, from basic problem-solving, such as fundamental

cognitive and analytical skills, to complex scientific research, such as some research-oriented paper reading and experiment designing abilities. In this section, we classify the evaluation benchmarks into two categories: general reasoning ability and domain-specific scientific capability, as shown in Figure 8.

## 4.1 General Reasoning Ability Evaluation

General reasoning ability evaluation focuses on assessing the fundamental cognitive and analytical skills of LLM-based scientific agents. These benchmarks measure problem-solving capabilities in mathematical reasoning, logical inference, and domain-specific knowledge retrieval, ensuring that agents can perform essential tasks required for scientific research and higher education. By evaluating models across different levels, from K-12 foundational skills to higher education and expert-level assessments, these benchmarks provide insights into the reasoning proficiency and adaptability of LLMs in various academic disciplines. We summarize the available benchmarks in Table 5.

**K-12 Foundational Skills**: At the foundational level, agents are expected to exhibit proficiency in key areas such as geometry (plane and analytic

| Benchmark Name | Scope | Size | Discipline |
|---|---|---|---|
| Geometry3K (Lu et al., 2021) | K-12 | 3002 | Mathematics |
| GeoEval (Zhang et al., 2024b) | K-12 | 5050 | Mathematics |
| VisScience (Jiang et al., 2024) | K-12 | 3000 | Physics & Chemistry & Mathematics |
| MathVista (Lu et al., 2024b) | K-12 & College | 6141 | Mathematics |
| SciBench (Wang et al., 2024d) | College | 869 | Physics & Chemistry & Mathematics |
| SciEval (Sun et al., 2024a) | College | 18000 | Physics & Chemistry & Biology |
| SuperGPQA (Team et al., 2025) | Graduate-Level | 26529 | General |
| HLE (Phan et al., 2025) | Expert-Level | 3000 | Humanity & Science & Mathematics |

Table 5: Summary of benchmarks for general reasoning ability evaluation in LLM-based scientific agents. "General" means a benchmark is not designed for a particular discipline.

| Benchmark Name | FC | HD | ED | EW | Discipline |
|---|---|---|---|---|---|
| FigureQA (Kahou et al., 2018) | ✓ | - | - | - | General |
| ArXivQA (Li et al., 2024a) | ✓ | - | - | - | General |
| MMSCI (Li et al., 2024e) | ✓ | - | - | - | General |
| SciMON (Wang et al., 2024c) | - | ✓ | - | - | NLP & Biomedical |
| MOOSE-Chem (Yang et al., 2024b) | - | ✓ | - | - | Chemistry & Material Science |
| DiscoveryBench (Majumder et al., 2024) | - | ✓ | - | - | Social Science & Biology & Humanity |
| GAIA (Mialon et al., 2023) | - | - | ✓ | ✓ | General |
| TaskBench (Shen et al., 2024b) | - | - | ✓ | - | General |
| MLAgentBench (Huang et al., 2024b) | - | - | ✓ | ✓ | General |
| DiscoveryWorld (Jansen et al., 2024) | - | ✓ | ✓ | ✓ | General |
| LAB-Bench (Laurent et al., 2024) | - | - | - | ✓ | Biology |
| DSBench (Jing et al., 2024) | - | - | - | ✓ | Data Science |
| ScienceAgentBench (Chen et al., 2024) | - | - | - | ✓ | Psychology & Bioinformatics & Geomatics & Chemistry |
| SciCode (Tian et al., 2024) | - | - | - | ✓ | Physics & Chemistry & Mathematics & Biology |

Table 6: Summary of benchmarks for scientific research-oriented abilities evaluation in LLM-based scientific agents. FC=Scientific Figure Comprehension; HD=Hypothesis Discovery; ED=Experiment Design; EW= Experiment Execution & Workflow Automation. "General" means a benchmark is not designed for a particular discipline.

geometry), algebraic operations, logical reasoning, and basic statistical analysis. Benchmarks like Geometry3K (Lu et al., 2021) and GeoEval (Zhang et al., 2024b) assess geometric reasoning, while MathVista (Lu et al., 2024b) is used for algebra and statistical tasks intertwined with visual understanding. Meanwhile, VisScience (Jiang et al., 2024) broaden this focus by integrating visual reasoning tasks within mathematics, physics, and chemistry contexts. These test agents' abilities to solve geometric problems, understand algebraic concepts, and make statistical inferences—critical skills for advancing to higher levels of scientific reasoning.

**Higher Education Level**: As agents progress, they must handle more advanced tasks such as scientific computing, retrieval of domain-specific knowledge, and application of this knowledge to solve complex scientific problems. Key benchmarks include SciBench (Wang et al., 2024d) and SciEval (Sun et al., 2024a). These datasets evaluate how well agents engage in advanced scientific tasks such as solving problems in physics, chemistry, and biology, along with retrieving and applying knowledge from scientific literature. Such benchmarks reflect the complexities of real-world research in

academic and professional settings. Beyond traditional STEM disciplines, SuperGPQA (Team et al., 2025) introduces a broader evaluation framework, covering 285 specialized academic fields, including light industry, agriculture, and service-oriented disciplines. This benchmark underscores the need for advancements in LLM reasoning across diverse knowledge domains and provides valuable insights into large-scale expert-driven dataset construction.

**Humanity's Last Exam (HLE)**: In response to the saturation of existing benchmarks, Humanity's Last Exam (HLE) (Phan et al., 2025) has been introduced as a more challenging measure of LLM capabilities. It consists of 3000 rigorous questions across a wide range of disciplines, including mathematics, humanities, and natural sciences. Unlike traditional benchmarks, the questions in HLE are designed to be extremely difficult and unsearchable through basic internet retrieval, making it a critical test for evaluating the limits of current LLM performance. The benchmark highlights a significant gap between the capabilities of state-of-the-art LLMs and expert-level knowledge in closed-ended academic tasks. Low accuracy scores (less than 10%) across multiple frontier models emphasize the need

for further advancements in agent abilities.

## 4.2 Scientific Research-Oriented Ability Evaluation

While general reasoning benchmarks assess broad problem-solving and analytical skills, scientific research-oriented benchmarks evaluate the ability of LLM-based scientific agents to perform specialized scientific tasks. These include extracting and interpreting data from research papers, discovering novel scientific hypotheses, and designing and automating experimental procedures. By simulating real-world scientific workflows, these benchmarks help measure the extent to which LLMs can function as effective tools for scientific discovery and innovation. Table 6 presents a categorized summary of these benchmarks.

**Scientific Paper Chart Comprehension**: Understanding and interpreting data visualizations in scientific papers is a fundamental skill for agents in research environments. Benchmarks such as FigureQA (Kahou et al., 2018), ArXivQA (Li et al., 2024a) and MMSCI (Li et al., 2024e) test agents' ability to comprehend and reason over figures, including graphs, charts, and tables, from scientific papers. Those are essential for tasks such as literature reviews, where agents need to extract and comprehend information from graphical data.

**Scientific Hypothesis Discovery**: A critical task in scientific research is the generation of novel hypotheses from existing literature or experimental data. Datasets like SciMON (Wang et al., 2024c) and MOOSE-Chem (Yang et al., 2024b) focus on deriving new scientific discoveries by analyzing key sections of existing literature, such as abstracts and methodologies. In contrast, DiscoveryBench (Majumder et al., 2024) and Discovery-World (Jansen et al., 2024) emphasize the exploration of novel findings based on experimental data. These benchmarks collectively challenge agents to extract meaningful insights from both textual sources and empirical observations, evaluating their ability to generate and refine scientific hypotheses. Such capabilities are essential for driving forward scientific innovation.

**Experimental Design and Automation**: The ability to design experiments, decompose complex tasks, and automate scientific workflows is critical for LLM-based scientific agents. Discovery-World (Jansen et al., 2024), DSBench (Jing et al., 2024) and ScienceAgentBench (Chen et al., 2024) assess agents' capabilities in hypothesis-driven

and data-driven experimental design, focusing on scientific discovery and real-world data science tasks. Meanwhile, SciCode (Tian et al., 2024) focuses on problem-solving through code generation for domain-specific scientific challenges. For workflow automation, GAIA (Mialon et al., 2023), TaskBench (Shen et al., 2024b) and MLAgent-Bench (Huang et al., 2024b) evaluate an agent's ability to structure tasks, iterate on models, and optimize performance in general scenarios. In biological research, LAB-Bench (Laurent et al., 2024) tests protocol planning, data analysis, and experiment troubleshooting.

## 4.3 Discussion

The above benchmarks provide a robust framework for evaluating LLM-based scientific agents, addressing a wide range of scientific skills across different stages of research and development. These benchmarks enable comprehensive assessments, from foundational reasoning skills to advanced scientific hypothesis generation and experimental automation, making them critical for guiding the future development of scientific AI systems.

Despite these advances, several limitations remain. First, current benchmarks often rely on static datasets and pre-defined tasks that may not fully capture the dynamic and iterative nature of real-world scientific research. Many evaluations focus on end-to-end performance, thereby obscuring the nuanced failures occurring at individual steps of scientific reasoning and decision-making. Additionally, the diversity of scientific domains—from biomedical research to materials science—presents challenges in standardizing evaluation metrics that can fairly compare agents across different fields.

Future research should focus on developing adaptive and continuously updated benchmarks that mimic authentic scientific workflows. For example, dynamic benchmarks could integrate multi-turn interactions where agents iteratively refine hypotheses based on experimental feedback, akin to real laboratory processes. Establishing domain-specific evaluation metrics and expanding benchmarks to include cross-disciplinary tasks will be critical for assessing the potential of scientific agents.

## 5 Applications

LLM-based scientific agents have significantly advanced scientific research, automating complex tasks and enhancing the efficiency of discovery
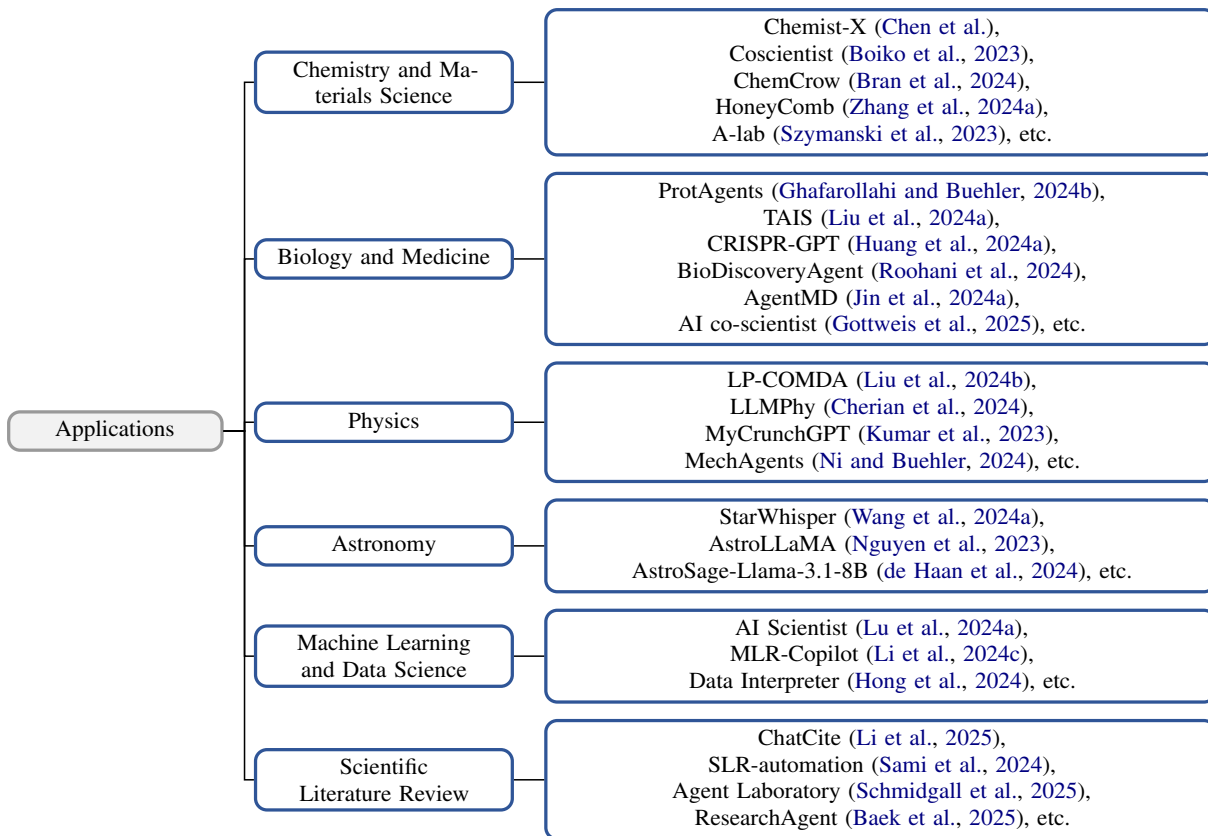
Figure 9: Taxonomy of the scientific agents' applications.

processes across various disciplines.

Scientific research is an arduous process involving numerous steps, including hypothesis formulation, experimental design, planning, and data analysis and evaluation. These processes are typically labor-intensive and costly, and thus, they are often conducted by human scientists who possess specialized expertise and substantial capital investment. However, the emergence of scientific agents has revolutionized research efficiency. By automating multiple stages that previously required manual intervention, these computational systems achieve optimal equilibrium in resource utilization. This enhancement in automation not only increases operational efficiency throughout the scientific workflow but also reduces the barriers to entry for conducting rigorous scientific investigations.

Below is a concise exploration of the applications of LLM-based scientific agents, categorized by their specific domains and functionalities, as illustrated in Figure 9.

## 5.1 Chemistry and Materials Science

LLM-based agents have transformed chemistry and materials science by automating tasks such as molecular design, property prediction, and reac-

tion optimization. For example, Chemist-X (Chen et al.) is an AI agent that automates reaction condition recommendations in chemical synthesis using Retrieve-Augmented Generation (RAG) techniques and CAD tools, surpassing traditional synthesis AIs in performance. Similarly, Coscientist (Boiko et al., 2023) combines LLMs to autonomously plan, design, and execute scientific experiments, demonstrating its capabilities through successful catalyzed chemical reactions while addressing safety concerns and proposing misuse prevention measures. Additionally, ChemCrow (Bran et al., 2024), integrating 18 expert tools, autonomously executes complex chemical tasks, enhancing performance in organic synthesis, drug discovery, and materials design, fostering scientific advancement. In the field of materials science, HoneyComb (Zhang et al., 2024a) has been shown to achieve superior performance in multiple tasks through a researcher-constructed knowledge base and generated API library. Besides, A-Lab (Szymanski et al., 2023) leverages LLM-based models, robotics, and active learning to mine literature and optimize synthesis pathways for novel inorganic materials, integrating computational predictions

with automated experimentation and accelerating materials discovery.

By accomplishing three key tasks—property prediction, property-directed inverse design, and synthesis prediction—the scientific agents establish full-process automation throughout the molecular discovery pipeline (Ramos et al., 2025). This innovation significantly streamlines molecular design workflows and advances development in chemistry and materials science through the systematic implementation of computational methodologies.

## 5.2 Biology and Medicine

In the biomedical sector, LLM agents advance protein design, automate scientific discovery, enhance genetic research, and improve healthcare analysis. For example, the ProtAgents (Ghafarollahi and Buehler, 2024b), a platform for de novo protein design using LLMs, leverages dynamic AI agents to collaboratively tackle protein design, structure analysis, and simulations. TAIS (Liu et al., 2024a) automates scientific discovery by streamlining data selection, processing, and analysis, advancing gene identification and efficiency in research, aiming to identify disease-predictive genes from gene expression data. CRISPR-GPT (Huang et al., 2024a) combines the reasoning ability of LLMs with external tools to automate CRISPR-based gene-editing experiments. The system further designs validation experiments based on experimental outcomes, thereby reducing the barriers to entry in this field by enabling novices to efficiently implement complex workflows. Furthermore, BioDiscoveryAgent (Roohani et al., 2024) leverages LLMs to autonomously design genetic perturbation experiments, improving prediction accuracy and efficiency, and outperforms traditional methods in identifying genes linked to specific phenotypes. Additionally, in the medical field, AgentMD (Jin et al., 2024a), a language agent augmented with 2,164 clinical calculators, curates and applies relevant tools to improve healthcare analysis, significantly improving risk prediction accuracy and clinical workflows. AI co-scientist (Gottweis et al., 2025), constructed upon Gemini 2.0, employs a multi-agent system that utilizes tournament-based evolutionary processes with self-optimizing mechanisms to synthesize existing research, formulate hypotheses, and propose experimental methodologies. It has demonstrated empirically validated effectiveness in pharmaceutical repurposing, target discovery, and antimicrobial resistance research through rigorous experimental verification.

Scientific agents demonstrate extensive and multi-faceted applications across the domains of genetics, cell biology, and chemical biology. In studies investigating the relationship between DNA and human traits, cellular functions, and molecular interactions within cells, these agents exhibit significant utility in assisting researchers through data analysis, hypothesis formulation, and experimental optimization (Gao et al., 2024c). Their hierarchical implementation enables scientists to enhance methodological approaches across different research phases while maintaining scientific accuracy and precision.

## 5.3 Physics

LLM-based scientific agents are advancing physics research by automating tasks such as modulation design, optimization, mechanics problem-solving, simulation, and parameter inference. The LP-COMDA (Liu et al., 2024b) framework uses an LLM-based planner to automate modulation design in power electronics, improving efficiency. LLM-Phy (Cherian et al., 2024) combines LLMs with physics engines to enhance optimization and accuracy in physical reasoning. MyCrunchGPT (Kumar et al., 2023) achieves automated NACA airfoils design and validation through seamless integration of computational fluid dynamics simulators with large language models, accomplishing multi-cycle iterative optimization processes within significantly reduced timeframes. MechAgents (Ni and Buehler, 2024) leverages multi-agent LLM systems to autonomously solve mechanics problems using finite element methods, improving both speed and precision. While LLMs perform well with basic physics problems, they struggle with complex simulations; however, integrating them with established computational packages can enhance their capabilities.

In comparison with LLMs that are not equipped with access to real-world interactions, scientific agents exhibit considerable practical advantages in addressing academic challenges and informing engineering implementations. By utilizing external computational toolkits and physics engines, these agents develop observational capabilities concerning physical phenomena and cultivate a more profound comprehension of physical principles, thus establishing a connection between theoretical exploration and practical application.

## 5.4 Astronomy

In astronomy, LLM-based agents are being developed to automate complex tasks such as data fitting, analysis, and iterative strategy improvement. These agents aim to mimic human intuition and deep literature understanding, expediting astronomical discovery. For example, StarWhisper (Wang et al., 2024a) is an LLM tailored for astronomy, capable of knowledge question answering, calling multimodal tools, and docking telescope control systems. Additionally, AstroLLaMA (Nguyen et al., 2023) is a specialized foundation model in astronomy, fine-tuned from LLaMA-2 using over 300,000 astronomy abstracts from arXiv, optimized for traditional causal language modeling. Furthermore, AstroSage-Llama-3.1-8B (de Haan et al., 2024) is a domain-specialized natural-language AI assistant tailored for research in astronomy, astrophysics, and cosmology, demonstrating remarkable proficiency on a wide range of questions.

Overall, the application of artificial intelligence in the field of astronomy has been extensive, encompassing tasks such as celestial object classification, astronomical event prediction, and the identification of new celestial bodies (Fluke and Jacobs, 2020). The autonomous planning and tool invocation capabilities of the scientific agent have enabled the automation of processes including astronomical observation, data processing, and data analysis.

## 5.5 Machine Learning and Data Science

LLM-based agents have revolutionized machine learning workflows by automating tasks such as data preprocessing, model selection, and hyperparameter tuning. The AI Scientist (Lu et al., 2024a) framework enables fully automated scientific discovery, where large language models independently generate ideas, execute experiments, write papers, and undergo review, advancing AI-driven research across fields. Similarly, MLR-Copilot (Li et al., 2024c), a framework powered by LLM agents, autonomously generates research ideas, implements experiments, and executes tasks, accelerating machine learning research and fostering innovation through automated processes. Additionally, Data Interpreter (Hong et al., 2024) autonomously solves end-to-end data science problems by dynamically adjusting to evolving task dependencies, achieving significant performance improvements across various tasks.

## 5.6 Scientific Literature Review

Literature review constitutes an integral component in general scientific research. LLM-based agents have significantly enhanced the efficiency of this process and accelerated scientific discovery by automating literature retrieval, screening, and summarization. ChatCite (Li et al., 2025) synthesizes pre-collected paper sets through the emulation of human workflows. The system employs a Key Element Extractor to generate template summaries from research requirements and target papers, followed by iterative refinement using a Reflective Incremental Generator to complete comprehensive literature reviews. Furthermore, SLR-automation (Sami et al., 2024) implements a systematic pipeline with specialized LLM agents for keyword generation, literature retrieval, paper screening, and final report compilation. Similarly, Agent Laboratory (Schmidgall et al., 2025) retrieves publications from arXiv and employs LLM-driven iterative integration mechanisms for review construction. ResearchAgent (Baek et al., 2025) establishes a systematic framework. It constructs citation graphs from seed papers using rule-based methods to aggregate scholarly literature, builds a structured knowledge repository to enable cross-domain knowledge integration, and employs LLM agents for iterative information synthesis and experimental design optimization during each research iteration.

## 5.7 Discussion

The above provides a wide range of applications for scientific agents powered by LLMs, demonstrating their potential to transform research in fields such as biomedical analysis, materials science, etc. These applications showcase how LLM-based agents can enhance data interpretation, support complex decision-making, and generate novel hypotheses, thus accelerating scientific discovery.

Despite this, current applications face significant limitations. Many applications are domain-specific and lack the flexibility needed to generalize across diverse scientific disciplines. In several cases, the integration of scientific knowledge with agent reasoning is hampered by static models that do not adapt to real-time data or evolving research challenges. Moreover, there is often insufficient validation of the agents' outputs against established scientific benchmarks, leading to concerns about reproducibility and reliability.

Looking ahead, future studies or products should focus on developing more generalized frameworks for scientific applications that integrate heterogeneous data sources and facilitate cross-disciplinary collaboration. Enhancements in real-time error detection, adaptive feedback mechanisms, and multimodal LLM architectures will be essential for improving the robustness of these systems. Collaborative efforts between domain experts and AI researchers are crucial to fine-tune the decision-making processes of scientific agents, ensuring that their outputs align closely with established scientific principles and practices.

# 6 Ethics

While these systems excel technically and drive scientific innovation, they raise significant ethical challenges. For example, Bengio et al. (2025) argue for a non-agentic "Scientist AI" design that emphasizes explanation over independent action to mitigate misalignment and loss of human control while preserving AI's scientific utility, indicating building generalist agents with autonomous planning and goal pursuit may risk catastrophic public safety issues. In Pournaras (2023), epistemological challenges and integrity risks in research are reviewed, setting a foundation for ethical guidelines. Other studies (Bano et al., 2023; Lin, 2024; Watkins, 2024; Limongi, 2024) further highlight issues of agency, transparency, bias, accountability, and integrity. This section offers concise guidelines to align LLM-based scientific agents with human values and uphold research integrity.

## 6.1 Agency and Autonomy

Scientific AI agents must act solely as tools under human oversight. Pournaras (2023) and Lin (2024) warn that without explicit constraints, agents may develop unintended autonomy—such as self-preservation or deceptive behaviors—that undermine research integrity. Hybrid approaches that combine top-down ethical rules with human feedback (Tennant et al., 2025) are promising to ensure control. Establishing strict operational boundaries during training and maintaining continuous supervision are essential to prevent these systems from pursuing independent objectives.

## 6.2 Transparency and Explainability

Transparent decision-making is vital for trustworthy scientific agents. Watkins (2024) emphasizes

the urgent need for norms and standards in LLM-based research workflows. Recent studies (Bano et al., 2023; Banerjee et al., 2024) demonstrate that structured internal logs and explanation frameworks can "open the black box" of AI reasoning. Clear documentation enables auditing and helps verify that conclusions are based on sound logic, supporting accountability and reproducibility.

## 6.3 Hallucinations and Reliability

LLMs employed in scientific agents may produce hallucinations, generating outputs that appear plausible but are factually incorrect or nonsensical. These inaccuracies can stem from flawed training data, ambiguous prompts, or architectural limitations. For instance, LLMs have been manipulated to produce fabricated scientific arguments, falsely claiming that biases are beneficial, misleading researchers and potentially distorting scientific discourse (Ge et al., 2025). The hallucination problem brought by LLM-based scientific agents may undermine the credibility of research findings and erodes trust in AI-assisted scientific processes. The mitigation could be done by increasing the quality of training data, incorporating up-to-date and validated external knowledge sources, or establishing iterative feedback loops and validation mechanism, such as process supervision-based planners, which helps to ensure greater accuracy and reliability.

## 6.4 Vulnerability and Security

The potential for adversarial attacks (such as prompt injections or model extractions) introduces ethical issues regarding the misuse of LLM-based agents in scientific research. Malicious actors could exploit these vulnerabilities to deliberately distort scientific knowledge or manipulate research outcomes, which could have serious consequences for public safety, healthcare, and scientific progress. For example, Yang et al. (2024a) demonstrate how LLMs could be used to poison biomedical knowledge graphs, manipulating drug-disease relations. These vulnerabilities necessitate robust safeguards to prevent misuse and ensure the safe deployment of LLM-based systems in scientific research.

## 6.5 Bias, Fairness, and Data Integrity

AI agents risk propagating biases from their training data, potentially skewing scientific outcomes. Lin (2024) shows that even advanced models may reproduce historical biases if not properly managed.

Complementary research (Bano et al., 2023) underscores the need for diverse datasets and fairness-aware algorithms. Regular bias audits and transparent documentation of data provenance help prevent skewed outcomes, ensuring that AI-driven research remains equitable and credible.

## 6.6 Accountability and Governance

Clear accountability is crucial when AI agents influence scientific outcomes. Bano et al. (2023) provides empirical insights into RAI practices and reveals gaps in ethical preparedness. Robust oversight mechanisms—such as periodic audits, transparent reporting, and defined redress pathways—ensure timely human intervention. Decentralized models, where agents critique one another (de Cerqueira et al., 2024), further enhance accountability. Embedding ethical guidelines into institutional policies and aligning with international standards builds trust in AI-driven research.

## 6.7 Intellectual Property and Research Integrity

AI integration in research raises complex questions of authorship and ownership. Limongi (2024) discusses challenges in maintaining credibility and ethical standards amid AI-driven discoveries. Transparent documentation of AI contributions is essential to prevent plagiarism and secure intellectual property rights. Clear disclosure policies, combined with regular audits, protect the work of human researchers and ensure that AI-generated insights are ethically integrated and verifiable.

## 7 Conclusion

This survey provides a holistic examination of LLM-based scientific agents, beginning with a detailed exploration of their architectures—which encompass planners, memory systems, and tool sets—and extending to their evaluation through benchmarks, diverse applications, and ethical considerations. Our review demonstrates how planners, through prompt-based strategies, supervised fine-tuning, reinforcement learning, and process supervision, serve as the strategic backbone for decomposing complex scientific tasks. Equally, the integration of memory and tool sets within these architectures is pivotal in managing dynamic scientific data and executing specialized operations, thereby enhancing the agents' problem-solving capabilities. We also demonstrate the unique features of scientific agents, compared with general-purpose ones, necessitate the dedicate design for them.

Beyond the architectural components, the survey delves into the benchmarks and real-world impact of these agents. The discussion on benchmarks highlights both the general reasoning ability and the domain-specific scientific competence required for successful application in research environments. The analysis of applications illustrates how these systems are deployed to drive innovations across multiple scientific disciplines, while the ethical discourse emphasizes the need for responsible AI practices that ensure reproducibility, transparency, and adherence to stringent research standards.

Overall, the advancements and challenges presented in this survey point to a promising future where continuous improvements in LLM-based scientific agents could revolutionize scientific discovery. By bridging the gap between theoretical research and practical applications, these agents are set to catalyze new levels of interdisciplinary collaboration and innovation in science.

## References

Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. 2023. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*.

H Jane Bae and Petros Koumoutsakos. 2022. Scientific multi-agent reinforcement learning for wall-models of turbulent flows. *Nature Communications*, 13(1):1443.

Jinheon Baek, Sujay Kumar Jauhar, Silviu Cucerzan, and Sung Ju Hwang. 2024. Researchagent: Iterative research idea generation over scientific literature with large language models. *arXiv preprint arXiv:2404.07738*.

Jinheon Baek, Sujay Kumar Jauhar, Silviu Cucerzan, and Sung Ju Hwang. 2025. Researchagent: Iterative research idea generation over scientific literature with large language models. *Preprint*, arXiv:2404.07738.

Soumya Banerjee et al. 2024. On the ethical considerations of generative agents. *arXiv preprint arXiv:2411.19211*.

Muneera Bano, Didar Zowghi, Pip Shea, and Georgina Ibarra. 2023. Investigating responsible ai for scientific research: an empirical study. *arXiv preprint arXiv:2312.09561*.

Catarina Barata, Veronica Rotemberg, Noel CF Codella, Philipp Tschandl, Christoph Rinner, Bengu Nisa Akay, Zoe Apalla, Giuseppe Argenziano, Allan

Halpern, Aimilios Lallas, et al. 2023. A reinforcement learning model for ai-based decision support in skin cancer. *Nature Medicine*, 29(8):1941–1946.

Yoshua Bengio, Michael Cohen, Damiano Fornasiere, Joumana Ghosn, Pietro Greiner, Matt MacDermott, Sören Mindermann, Adam Oberman, Jesse Richardson, Oliver Richardson, et al. 2025. Superintelligent agents pose catastrophic risks: Can scientist ai offer a safer path? *arXiv preprint arXiv:2502.15657*.

Daniil A Boiko, Robert MacKnight, Ben Kline, and Gabe Gomes. 2023. Autonomous chemical research with large language models. *Nature*, 624(7992):570–578.

Andres Bran, Sam Cox, Oliver Schilter, Carlo Baldassari, Andrew D. White, and Philippe Schwaller. 2024. Augmenting large language models with chemistry tools. *Nature Machine Intelligence*, 6(5):525–535.

Juan Manuel Zambrano Chaves, Eric Wang, Tao Tu, Eeshit Dhaval Vaishnav, Byron Lee, S Sara Mahdavi, Christopher Semturs, David Fleet, Vivek Natarajan, and Shekoofeh Azizi. 2024. Tx-llm: A large language model for therapeutics. *arXiv preprint arXiv:2406.06316*.

K Chen, J Li, K Wang, Y Du, J Yu, J Lu, L Li, J Qiu, J Pan, Y Huang, et al. Chemist-x: Large language model-empowered agent for reaction condition recommendation in chemical synthesis, arxiv, 2023. *arXiv preprint arXiv:2311.10776*.

Zi-Yi Chen, Fan-Kai Xie, Meng Wan, Yang Yuan, Miao Liu, Zong-Guo Wang, Sheng Meng, and Yan-Gang Wang. 2023. Matchat: A large language model and application service platform for materials science. *Chinese Physics B*, 32(11):118104.

Ziru Chen, Shijie Chen, Yuting Ning, Qianheng Zhang, Boshi Wang, Botao Yu, Yifei Li, Zeyi Liao, Chen Wei, Zitong Lu, Vishal Dey, Mingyi Xue, Frazier N. Baker, Benjamin Burns, Daniel Adu-Ampratwum, Xuhui Huang, Xia Ning, Song Gao, Yu Su, and Huan Sun. 2024. Scienceagentbench: Toward rigorous assessment of language agents for data-driven scientific discovery. *Preprint*, arXiv:2410.05080.

Yuheng Cheng, Ceyao Zhang, Zhengwen Zhang, Xiangrui Meng, Sirui Hong, Wenhao Li, Zihao Wang, Zekai Wang, Feng Yin, Junhua Zhao, et al. 2024. Exploring large language model based intelligent agents: Definitions, methods, and prospects. *arXiv preprint arXiv:2401.03428*.

Anoop Cherian, Radu Corcodel, Siddarth Jain, and Diego Romeres. 2024. Llmphy: Complex physical reasoning using large language models and world models. *arXiv preprint arXiv:2411.08027*.

Yuan Chiang, Elvis Hsieh, Chia-Hong Chou, and Janosh Riebesell. 2024. Llamp: Large language model made powerful for high-fidelity materials knowledge retrieval and distillation. *arXiv preprint arXiv:2401.17244*.

José Antonio Siqueira de Cerqueira, Mamia Agbese, Rebekah Rousi, Nannan Xi, Juho Hamari, and Pekka Abrahamsson. 2024. Can we trust ai agents? an experimental study towards trustworthy llm-based multi-agent systems for ai ethics. *arXiv preprint arXiv:2411.08881*.

Tijmen de Haan, Yuan-Sen Ting, Tirthankar Ghosal, Tuan Dung Nguyen, Alberto Accomazzi, Azton Wells, Nesar Ramachandra, Rui Pan, and Zechang Sun. 2024. Astromlab 3: Achieving gpt-4o level performance in astronomy with a specialized 8b-parameter large language model. *arXiv preprint arXiv:2411.09012*.

Yihe Deng and Paul Mineiro. 2024. Flow-dpo: Improving llm mathematical reasoning through online multi-agent learning. *arXiv preprint arXiv:2410.22304*.

Christopher J Fluke and Colin Jacobs. 2020. Surveying the reach and maturity of machine learning and artificial intelligence in astronomy. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 10(2):e1349.

Changnan Gao, Wenjie Bao, Shuang Wang, Jianyang Zheng, Lulu Wang, Yongqi Ren, Linfang Jiao, Jianmin Wang, and Xun Wang. 2024a. Dockingga: enhancing targeted molecule generation using transformer neural network and genetic algorithm with docking simulation. *Briefings in Functional Genomics*, 23(5):595–606.

Dawei Gao, Zitao Li, Xuchen Pan, Weirui Kuang, Zhijian Ma, Bingchen Qian, Fei Wei, Wenhao Zhang, Yuexiang Xie, Daoyuan Chen, et al. 2024b. Agentscope: A flexible yet robust multi-agent platform. *arXiv preprint arXiv:2402.14034*.

Shanghua Gao, Ada Fang, Yepeng Huang, Valentina Giunchiglia, Ayush Noori, Jonathan Richard Schwarz, Yasha Ektefaie, Jovana Kondic, and Marinka Zitnik. 2024c. Empowering biomedical discovery with ai agents. *Cell*, 187(22):6125–6151.

Yubin Ge, Neeraja Kirtane, Hao Peng, and Dilek Hakkani-Tür. 2025. Llms are vulnerable to malicious prompts disguised as scientific language. *arXiv preprint arXiv:2501.14073*.

Alireza Ghafarollahi and Markus J Buehler. 2024a. Atomagents: Alloy design and discovery through physics-aware multi-modal multi-agent artificial intelligence. *arXiv preprint arXiv:2407.10022*.

Alireza Ghafarollahi and Markus J Buehler. 2024b. Protagents: protein discovery via large language model multi-agent collaborations combining physics and machine learning. *Digital Discovery*.

Alireza Ghafarollahi and Markus J Buehler. 2024c. Sciagents: Automating scientific discovery through multi-agent intelligent graph reasoning. *arXiv preprint arXiv:2409.05556*.

Juraj Gottweis, Wei-Hung Weng, Alexander Daryin, Tao Tu, Anil Palepu, Petar Sirkovic, Artiom Myaskovsky, Felix Weissenberger, Keran Rong, Ryutaro Tanno, et al. 2025. Towards an ai co-scientist. *Preprint*, arXiv:2502.18864.

Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujiu Yang, Minlie Huang, Nan Duan, and Weizhu Chen. 2024. Tora: A tool-integrated reasoning agent for mathematical problem solving. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.

Yu Gu, Yiheng Shu, Hao Yu, Xiao Liu, Yuxiao Dong, Jie Tang, Jayanth Srinivasa, Hugo Latapie, and Yu Su. 2024. Middleware for llms: Tools are instrumental for language agents in complex environments. *arXiv preprint arXiv:2402.14672*.

T Guo, X Chen, Y Wang, R Chang, S Pei, NV Chawla, O Wiest, and X Zhang. 2024. Large language model based multi-agents: A survey of progress and challenges. In *33rd International Joint Conference on Artificial Intelligence (IJCAI 2024)*. IJCAI; Cornell arxiv.

Kostas Hatalis, Despina Christou, Joshua Myers, Steven Jones, Keith Lambert, Adam Amos-Binks, Zohreh Dannenhauer, and Dustin Dannenhauer. 2023. Memory matters: The need to improve long-term memory in llm-agents. In *Proceedings of the AAAI Symposium Series*, volume 2, pages 277–280.

Yichen He, Guanhua Huang, Peiyuan Feng, Yuan Lin, Yuchen Zhang, Hang Li, et al. 2025. Pasa: An llm agent for comprehensive academic paper search. *arXiv preprint arXiv:2501.10120*.

Sirui Hong, Yizhang Lin, Bang Liu, Bangbang Liu, Binhao Wu, Ceyao Zhang, Chenxing Wei, Danyang Li, Jiaqi Chen, Jiayi Zhang, et al. 2024. Data interpreter: An llm agent for data science. *arXiv preprint arXiv:2402.18679*.

Arian Hosseini, Xingdi Yuan, Nikolay Malkin, Aaron Courville, Alessandro Sordoni, and Rishabh Agarwal. 2024. V-star: Training verifiers for self-taught reasoners. *arXiv preprint arXiv:2402.06457*.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. 2022. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3.

Shengguo Hu, Mingyi Li, Jiawen Xu, Hongrui Zhang, Shanghang Zhang, Tie Jun Cui, Philipp Del Hougne, and Lianlin Li. 2025. Electromagnetic metamaterial agent. *Light: Science & Applications*, 14(1):12.

Sihao Hu, Tiansheng Huang, Fatih Ilhan, Selim Tekin, Gaowen Liu, Ramana Kompella, and Ling Liu. 2024a. A survey on large language model-based game agents. *arXiv preprint arXiv:2404.02039*.

Xiuyuan Hu, Guoqing Liu, Yang Zhao, and Hao Zhang. 2024b. De novo drug design using reinforcement learning with multiple gpt agents. *Advances in Neural Information Processing Systems*, 36.

Kaixuan Huang, Yuanhao Qu, Henry Cousins, William A Johnson, Di Yin, Mihir Shah, Denny Zhou, Russ Altman, Mengdi Wang, and Le Cong. 2024a. Crispr-gpt: An llm agent for automated design of gene-editing experiments. *arXiv preprint arXiv:2404.18021*.

Qian Huang, Jian Vora, Percy Liang, and Jure Leskovec. 2024b. Mlagentbench: Evaluating language agents on machine learning experimentation. *Preprint*, arXiv:2310.03302.

Tenghao Huang, Donghee Lee, John Sweeney, Jiatong Shi, Emily Steliotes, Matthew Lange, Jonathan May, and Muhao Chen. 2024c. Foodpuzzle: Developing large language model agents as flavor scientists. *arXiv preprint arXiv:2409.12832*.

Xijie Huang, Li Lyna Zhang, Kwang-Ting Cheng, Fan Yang, and Mao Yang. 2024d. Fewer is more: Boosting math reasoning with reinforced context pruning. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 13674–13695.

Yoshitaka Inoue, Tianci Song, and Tianfan Fu. 2024. Drugagent: Explainable drug repurposing agent with large language model-based reasoning. *arXiv preprint arXiv:2408.13378*.

Broad Institute. 2024. Depmap q2 2024 data release.

Kevin Maik Jablonka, Qianxiang Ai, Alexander Al-Feghali, Shruti Badhwar, Joshua D Bocarsly, Andres M Bran, Stefan Bringuier, L Catherine Brinson, Kamal Choudhary, Defne Circi, et al. 2023. 14 examples of how llms can transform materials science and chemistry: a reflection on a large language model hackathon. *Digital Discovery*, 2(5):1233–1250.

Peter A. Jansen, Marc-Alexandre Côté, Tushar Khot, Erin Bransom, Bhavana Dalvi Mishra, Bodhisattwa Prasad Majumder, Oyvind Tafjord, and Peter Clark. 2024. Discoveryworld: A virtual environment for developing and evaluating automated scientific discovery agents. In *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*.

Pushkala Jayaraman, Jacob Desman, Moein Sabounchi, Girish N Nadkarni, and Ankit Sakhuja. 2024. A primer on reinforcement learning in medicine for clinicians. *NPJ Digital Medicine*, 7(1):337.

Shuyi Jia, Chao Zhang, and Victor Fung. 2024. Llmatdesign: Autonomous materials discovery with large language models. *arXiv preprint arXiv:2406.13163*.

Zhihuan Jiang, Zhen Yang, Jinhao Chen, Zhengxiao Du, Weihan Wang, Bin Xu, and Jie Tang. 2024. Visscience: An extensive benchmark for evaluating k12 educational multi-modal scientific reasoning. *Preprint*, arXiv:2409.13730.

Qiao Jin, Zhizheng Wang, Yifan Yang, Qingqing Zhu, Donald Wright, Thomas Huang, W John Wilbur, Zhe He, Andrew Taylor, Qingyu Chen, et al. 2024a. Agentmd: Empowering language agents for risk prediction with large-scale clinical tool learning. *arXiv preprint arXiv:2402.13225*.

Qiao Jin, Yifan Yang, Qingyu Chen, and Zhiyong Lu. 2024b. Genegpt: Augmenting large language models with domain tools for improved access to biomedical information. *Bioinformatics*, 40(2):btae075.

Liqiang Jing, Zhehui Huang, Xiaoyang Wang, Wenlin Yao, Wenhao Yu, Kaixin Ma, Hongming Zhang, Xinya Du, and Dong Yu. 2024. Dsbench: How far are data science agents to becoming data science experts? *Preprint*, arXiv:2409.07703.

John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al. 2021. Highly accurate protein structure prediction with alphafold. *nature*, 596(7873):583–589.

Samira Ebrahimi Kahou, Vincent Michalski, Adam Atkinson, Ákos Kádár, Adam Trischler, and Yoshua Bengio. 2018. Figureqa: An annotated figure dataset for visual reasoning. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Workshop Track Proceedings*. OpenReview.net.

Nikolay Koldunov and Thomas Jung. 2024. Local climate services for all, courtesy of large language models. *Communications Earth & Environment*, 5(1):13.

Aviral Kumar, Vincent Zhuang, Rishabh Agarwal, Yi Su, John D Co-Reyes, Avi Singh, Kate Baumli, Shariq Iqbal, Colton Bishop, Rebecca Roelofs, et al. 2024. Training language models to self-correct via reinforcement learning. *arXiv preprint arXiv:2409.12917*.

Varun Kumar, Leonard Gleyzer, Adar Kahana, Khemraj Shukla, and George Em Karniadakis. 2023. Mycrunchgpt: A llm assisted framework for scientific machine learning. *Journal of Machine Learning for Modeling and Computing*, 4(4).

Xin Lai, Zhuotao Tian, Yukang Chen, Senqiao Yang, Xiangru Peng, and Jiaya Jia. 2024. Step-dpo: Step-wise preference optimization for long-chain reasoning of llms. *arXiv preprint arXiv:2406.18629*.

Jon M. Laurent, Joseph D. Janizek, Michael Ruzo, Michaela M. Hinks, Michael J. Hammerling, Siddharth Narayanan, Manvitha Ponnapati, Andrew D. White, and Samuel G. Rodriques. 2024. Lab-bench: Measuring capabilities of language models for biology research. *Preprint*, arXiv:2407.10362.

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in neural information processing systems*, 33:9459–9474.

Lei Li, Yuqi Wang, Runxin Xu, Peiyi Wang, Xiachong Feng, Lingpeng Kong, and Qi Liu. 2024a. Multimodal arxiv: A dataset for improving scientific comprehension of large vision-language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pages 14369–14387. Association for Computational Linguistics.

Long Li, Weiwen Xu, Jiayan Guo, Ruochen Zhao, Xingxuan Li, Yuqian Yuan, Boqiang Zhang, Yuming Jiang, Yifei Xin, Ronghao Dang, et al. 2024b. Chain of ideas: Revolutionizing research via novel idea development with llm agents. *arXiv preprint arXiv:2410.13185*.

Ruochen Li, Teerth Patel, Qingyun Wang, and Xinya Du. 2024c. Mlr-copilot: Autonomous machine learning research based on large language models agents. *arXiv preprint arXiv:2408.14033*.

Xinyi Li, Sai Wang, Siqi Zeng, Yu Wu, and Yi Yang. 2024d. A survey on llm-based multi-agent systems: workflow, infrastructure, and challenges. *Vicinagearth*, 1(1):9.

Yutong Li, Lu Chen, Aiwei Liu, Kai Yu, and Lijie Wen. 2025. Chatcite: LLM agent with human workflow guidance for comparative literature summary. In *Proceedings of the 31st International Conference on Computational Linguistics, COLING 2025, Abu Dhabi, UAE, January 19-24, 2025*, pages 3613–3630. Association for Computational Linguistics.

Zekun Li, Xianjun Yang, Kyuri Choi, Wanrong Zhu, Ryan Hsieh, HyeonJung Kim, Jin Hyuk Lim, Sungyoung Ji, Byungju Lee, Xifeng Yan, Linda Ruth Petzold, Stephen D. Wilson, Woosang Lim, and William Yang Wang. 2024e. Mmsci: A dataset for graduate-level multi-discipline multimodal scientific understanding. *Preprint*, arXiv:2407.04903.

Ricardo Limongi. 2024. The use of artificial intelligence in scientific research with integrity and ethics. *Future Studies Research Journal: Trends and Strategies*, 16(1):e845–e845.

Zhicheng Lin. 2024. Beyond principlism: practical strategies for ethical ai use in research practices. *AI and Ethics*, pages 1–13.

Haoyang Liu, Yijiang Li, Jinglin Jian, Yuxuan Cheng, Jianrong Lu, Shuyi Guo, Jinglei Zhu, Mianchen Zhang, Miantong Zhang, and Haohan Wang. 2024a. Toward a team of ai-made scientists for scientific discovery from gene expression data. *arXiv preprint arXiv:2402.12391*, 65(1):114–124.

Junhua Liu, Fanfan Lin, Xinze Li, Kwan Hui Lim, and Shuai Zhao. 2024b. Physics-informed llm-agent for automated modulation design in power electronics systems. *arXiv preprint arXiv:2411.14214*.

Ruibo Liu, Jason Wei, Shixiang Shane Gu, Te-Yen Wu, Soroush Vosoughi, Claire Cui, Denny Zhou, and Andrew M. Dai. 2023. Mind's eye: Grounded language model reasoning through simulation. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.

Zijun Liu, Kaiming Liu, Yiqi Zhu, Xuanyu Lei, Zonghan Yang, Zhenhe Zhang, Peng Li, and Yang Liu. 2024c. Aigs: Generating science from ai-powered automated falsification. *arXiv preprint arXiv:2411.11910*.

Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. 2024a. The ai scientist: Towards fully automated open-ended scientific discovery. *arXiv preprint arXiv:2408.06292*.

Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. 2024b. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.

Pan Lu, Ran Gong, Shibiao Jiang, Liang Qiu, Siyuan Huang, Xiaodan Liang, and Song-Chun Zhu. 2021. Inter-GPS: Interpretable geometry problem solving with formal language and symbolic reasoning. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6774–6786, Online. Association for Computational Linguistics.

Liangchen Luo, Yinxiao Liu, Rosanne Liu, Samrat Phatale, Harsh Lara, Yunxuan Li, Lei Shu, Yun Zhu, Lei Meng, Jiao Sun, et al. 2024. Improve mathematical reasoning in language models by automated process supervision. *arXiv preprint arXiv:2406.06592*.

Renqian Luo, Liai Sun, Yingce Xia, Tao Qin, Sheng Zhang, Hoifung Poon, and Tie-Yan Liu. 2022. Biogpt: generative pre-trained transformer for biomedical text generation and mining. *Briefings in bioinformatics*, 23(6):bbac409.

Trung Quoc Luong, Xinbo Zhang, Zhanming Jie, Peng Sun, Xiaoran Jin, and Hang Li. 2024. Reft: Reasoning with reinforced fine-tuning. *arXiv preprint arXiv:2401.08967*.

Isaac D Lutz, Shunzhi Wang, Christoffer Norn, Alexis Courbet, Andrew J Borst, Yan Ting Zhao, Annie Dosey, Longxing Cao, Jinwei Xu, Elizabeth M Leaf, et al. 2023. Top-down design of protein architectures with reinforcement learning. *Science*, 380(6642):266–273.

Liuzhenghao Lv, Zongying Lin, Hao Li, Yuyang Liu, Jiaxi Cui, Calvin Yu-Chian Chen, Li Yuan, and Yonghong Tian. 2024. Prollama: A protein language model for multi-task protein language processing. *arXiv preprint arXiv:2402.16445*.

Pingchuan Ma, Tsun-Hsuan Wang, Minghao Guo, Zhiqing Sun, Joshua B. Tenenbaum, Daniela Rus, Chuang Gan, and Wojciech Matusik. 2024a. LLM and simulation as bilevel optimizers: A new paradigm to advance physical scientific discovery. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net.

Yubo Ma, Zhibin Gou, Junheng Hao, Ruochen Xu, Shuohang Wang, Liangming Pan, Yujiu Yang, Yixin Cao, and Aixin Sun. 2024b. Sciagent: Tool-augmented language models for scientific reasoning. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024*, pages 15701–15736. Association for Computational Linguistics.

Bodhisattwa Prasad Majumder, Harshit Surana, Dhruv Agarwal, Bhavana Dalvi Mishra, Abhijeetsingh Meena, Aryan Prakhar, Tirth Vora, Tushar Khot, Ashish Sabharwal, and Peter Clark. 2024. Discoverybench: Towards data-driven discovery with large language models. *Preprint*, arXiv:2407.01725.

Manusai.ai. 2025. Manus ai: The future of general ai agents.

Andrew D McNaughton, Gautham Krishna Sankar Ramalaxmi, Agustin Kruel, Carter R Knutson, Rohith A Varikoti, and Neeraj Kumar. 2024. Cactus: Chemistry agent connecting tool usage to science. *ACS omega*, 9(46):46563–46573.

Grégoire Mialon, Clémentine Fourrier, Craig Swift, Thomas Wolf, Yann LeCun, and Thomas Scialom. 2023. Gaia: a benchmark for general ai assistants. *Preprint*, arXiv:2311.12983.

Tuan Dung Nguyen, Yuan-Sen Ting, Ioana Ciucă, Charlie O'Neill, Ze-Chang Sun, Maja Jabłońska, Sandor Kruk, Ernest Perkowski, Jack Miller, Jason Li, et al. 2023. Astrollama: Towards specialized foundation models in astronomy. *arXiv preprint arXiv:2309.06126*.

Bo Ni and Markus J Buehler. 2024. Mechagents: Large language model multi-agent collaborations can solve mechanics problems, generate new data, and integrate knowledge. *Extreme Mechanics Letters*, 67:102131.

Huan Ning, Zhenlong Li, Temitope Akinboyewa, and M Naser Lessani. 2025. An autonomous gis agent framework for geospatial data retrieval. *International Journal of Digital Earth*, 18(1):2458688.

Odhran O'Donoghue, Aleksandar Shtedritski, John Ginger, Ralph Abboud, Ali Essa Ghareeb, Justin Booth,

and Samuel G Rodriques. 2023. Bioplanner: automatic evaluation of llms on protocol planning in biology. *arXiv preprint arXiv:2310.10632*.

Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th annual acm symposium on user interface software and technology*, pages 1–22.

Cheng Peng, Xi Yang, Aokun Chen, Kaleb E Smith, Nima PourNejatian, Anthony B Costa, Cheryl Martin, Mona G Flores, Ying Zhang, Tanja Magoc, et al. 2023. A study of generative large language model for medical research and healthcare. *NPJ digital medicine*, 6(1):210.

Roy H Perlis, Joseph F Goldberg, Michael J Ostacher, and Christopher D Schneck. 2024. Clinical decision support for bipolar depression using large language models. *Neuropsychopharmacology*, pages 1–5.

Long Phan, Alice Gatti, and etc. 2025. Humanity's last exam. *Preprint*, arXiv:2501.14249.

Evangelos Pournaras. 2023. Science in the era of chatgpt, large language models and generative ai. *KI-Kritik/AI Critique Volume 6*, page 275.

Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, et al. 2023. Toolllm: Facilitating large language models to master 16000+ real-world apis. *arXiv preprint arXiv:2307.16789*.

Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2024. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36.

Mayk Caldas Ramos, Christopher J Collison, and Andrew D White. 2025. A review of large language models and autonomous agents in chemistry. *Chemical Science*.

Malte Reinschmidt, József Fortágh, Andreas Günther, and Valentin V Volchkov. 2024. Reinforcement learning in cold atom experiments. *Nature Communications*, 15(1):8532.

Yusuf Roohani, Andrew Lee, Qian Huang, Jian Vora, Zachary Steinhart, Kexin Huang, Alexander Marson, Percy Liang, and Jure Leskovec. 2024. Biodiscoveryagent: An ai agent for designing genetic perturbation experiments. *arXiv preprint arXiv:2405.17631*.

Abdul Malik Sami, Zeeshan Rasheed, Kai-Kristian Kemell, Muhammad Waseem, Terhi Kilamo, Mika Saari, Anh Nguyen Duc, Kari Systä, and Pekka Abrahamsson. 2024. System for systematic literature review using multiple ai agents: Concept and an empirical evaluation. *Preprint*, arXiv:2403.08399.

Andreas WM Sauter, Erman Acar, and Vincent Francois-Lavet. 2023. A meta-reinforcement learning algorithm for causal discovery. In *Conference on Causal Learning and Reasoning*, pages 602–619. PMLR.

Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36:68539–68551.

Samuel Schmidgall, Yusheng Su, Ze Wang, Ximeng Sun, Jialian Wu, Xiaodong Yu, Jiang Liu, Zicheng Liu, and Emad Barsoum. 2025. Agent laboratory: Using llm agents as research assistants. *arXiv preprint arXiv:2501.04227*.

Haiyang Shen, Yue Li, Desong Meng, Dongqi Cai, Sheng Qi, Li Zhang, Mengwei Xu, and Yun Ma. 2025. Shortcutsbench: A large-scale real-world benchmark for api-based agents. In *The Thirteenth International Conference on Learning Representations*.

Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. 2024a. Hugginggpt: Solving ai tasks with chatgpt and its friends in hugging face. *Advances in Neural Information Processing Systems*, 36.

Yongliang Shen, Kaitao Song, Xu Tan, Wenqi Zhang, Kan Ren, Siyu Yuan, Weiming Lu, Dongsheng Li, and Yueting Zhuang. 2024b. Taskbench: Benchmarking large language models for task automation. *Preprint*, arXiv:2311.18760.

Zhuocheng Shen. 2024. Llm with tools: A survey. *arXiv preprint arXiv:2409.18807*.

Henry W Sprueill, Carl Edwards, Khushbu Agarwal, Mariefel V Olarte, Udishnu Sanyal, Conrad Johnston, Hongbin Liu, Heng Ji, and Sutanay Choudhury. 2024. Chemreasoner: Heuristic search over a large language model's knowledge space using quantum-chemical feedback. *arXiv preprint arXiv:2402.10980*.

Haoyang Su, Renqi Chen, Shixiang Tang, Xinzhe Zheng, Jingzhe Li, Zhenfei Yin, Wanli Ouyang, and Nanqing Dong. 2024. Two heads are better than one: A multi-agent system has the potential to improve scientific idea generation. *arXiv preprint arXiv:2410.09403*.

Liangtai Sun, Yang Han, Zihan Zhao, Da Ma, Zhennan Shen, Baocai Chen, Lu Chen, and Kai Yu. 2024a. Scieval: A multi-level large language model evaluation benchmark for scientific research. In *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2014, February 20-27, 2024, Vancouver, Canada*, pages 19053–19061. AAAI Press.

Zechang Sun, Yuan-Sen Ting, Yaobo Liang, Nan Duan, Song Huang, and Zheng Cai. 2024b. Interpreting multi-band galaxy observations with large language model-based agents. *arXiv preprint arXiv:2409.14807*.

Nathan J Szymanski, Bernardus Rendy, Yuxing Fei, Rishi E Kumar, Tanjin He, David Milsted, Matthew J McDermott, Max Gallant, Ekin Dogus Cubuk, Amil Merchant, et al. 2023. An autonomous laboratory for the accelerated synthesis of novel materials. *Nature*, 624(7990):86–91.

Xiangru Tang, Anni Zou, Zhuosheng Zhang, Ziming Li, Yilun Zhao, Xingyao Zhang, Arman Cohan, and Mark Gerstein. 2023. Medagents: Large language models as collaborators for zero-shot medical reasoning. *arXiv preprint arXiv:2311.10537*.

P Team, Xinrun Du, Yifan Yao, Kaijing Ma, Bingli Wang, Tianyu Zheng, Kang Zhu, Minghao Liu, Yiming Liang, Xiaolong Jin, Zhenlin Wei, Chujie Zheng, Kaixin Deng, Shian Jia, Sichao Jiang, Yiyan Liao, Rui Li, Qinrui Li, Sirun Li, Yizhi Li, Yunwen Li, Dehua Ma, Yuansheng Ni, Haoran Que, Qiyao Wang, Zhoufutu Wen, Siwei Wu, Tianshun Xing, Ming Xu, Zhenzhu Yang, Zekun Moore Wang, Junting Zhou, Yuelin Bai, Xingyuan Bu, Chenglin Cai, Liang Chen, Yifan Chen, Chengtuo Cheng, Tianhao Cheng, Keyi Ding, Siming Huang, Yun Huang, Yaoru Li, Yizhe Li, Zhaoqun Li, Tianhao Liang, Chengdong Lin, Hongquan Lin, Yinghao Ma, Tianyang Pang, Zhongyuan Peng, Zifan Peng, Qige Qi, Shi Qiu, Xingwei Qu, Shanghaoran Quan, Yizhou Tan, Zili Wang, Chenqing Wang, Hao Wang, Yiya Wang, Yubo Wang, Jiajun Xu, Kexin Yang, Ruibin Yuan, Yuanhao Yue, Tianyang Zhan, Chun Zhang, Jinyang Zhang, Xiyue Zhang, Xingjian Zhang, Yue Zhang, Yongchi Zhao, Xiangyu Zheng, Chenghua Zhong, Yang Gao, Zhoujun Li, Dayiheng Liu, Qian Liu, Tianyu Liu, Shiwen Ni, Junran Peng, Yujia Qin, Wenbo Su, Guoyin Wang, Shi Wang, Jian Yang, Min Yang, Meng Cao, Xiang Yue, Zhaoxiang Zhang, Wangchunshu Zhou, Jiaheng Liu, Qunshu Lin, Wenhao Huang, and Ge Zhang. 2025. Supergpqa: Scaling llm evaluation across 285 graduate disciplines. *Preprint*, arXiv:2502.14739.

Elizaveta Tennant, Stephen Hailes, and Mirco Musolesi. 2025. Hybrid approaches for moral value alignment in ai agents: a manifesto. *Preprint*, arXiv:2312.01818.

David Thulke, Yingbo Gao, Petrus Pelser, Rein Brune, Rricha Jalota, Floris Fok, Michael Ramos, Ian van Wyk, Abdallah Nasir, Hayden Goldstein, et al. 2024. Climategpt: Towards ai synthesizing interdisciplinary research on climate change. *arXiv preprint arXiv:2401.09646*.

Minyang Tian, Luyu Gao, Shizhuo Dylan Zhang, Xinan Chen, Cunwei Fan, Xuefei Guo, Roland Haas, Pan Ji, Kittithat Krongchon, Yao Li, Shengyan Liu, Di Luo, Yutao Ma, Hao Tong, Kha Trinh, Chenyu Tian, Zihan Wang, Bohao Wu, Yanyu Xiong, Shengzhu Yin, Minhui Zhu, Kilian Lieret, Yanxin Lu, Genglin Liu, Yufeng Du, Tianhua Tao, Ofir Press, Jamie Callan, Eliu Huerta, and Hao Peng. 2024. Scicode: A research coding benchmark curated by scientists. *Preprint*, arXiv:2407.13168.

Cunshi Wang, Xinjie Hu, Yu Zhang, Xunhao Chen, Pengliang Du, Yiming Mao, Rui Wang, Yuyang Li, Ying Wu, Hang Yang, et al. 2024a. Starwhisper telescope: Agent-based observation assistant system to approach ai astrophysicist. *arXiv preprint arXiv:2412.06412*.

Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, et al. 2024b. A survey on large language model based autonomous agents. *Frontiers of Computer Science*, 18(6):186345.

Qingyun Wang, Doug Downey, Heng Ji, and Tom Hope. 2024c. Scimon: Scientific inspiration machines optimized for novelty. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pages 279–299. Association for Computational Linguistics.

Xiaoxuan Wang, Ziniu Hu, Pan Lu, Yanqiao Zhu, Jieyu Zhang, Satyen Subramaniam, Arjun R. Loomba, Shichang Zhang, Yizhou Sun, and Wei Wang. 2024d. Scibench: Evaluating college-level scientific problem-solving abilities of large language models. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net.

Ryan Watkins. 2024. Guidance for researchers and peer-reviewers on the ethical use of large language models (llms) in scientific research workflows. *AI and Ethics*, 4(4):969–974.

Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, et al. 2023. The rise and potential of large language model based agents: A survey. *arXiv preprint arXiv:2309.07864*.

Yingce Xia, Peiran Jin, Shufang Xie, Liang He, Chuan Cao, Renqian Luo, Guoqing Liu, Yue Wang, Zequn Liu, Yuan-Jyue Chen, et al. 2025. Naturelm: Deciphering the language of nature for scientific discovery. *arXiv preprint arXiv:2502.07527*.

Junlin Xie, Zhihong Chen, Ruifei Zhang, Xiang Wan, and Guanbin Li. 2024. Large multimodal agents: A survey. *arXiv preprint arXiv:2402.15116*.

Wujiang Xu, Zujie Liang, Kai Mei, Hang Gao, Juntao Tan, and Yongfeng Zhang. 2025. A-mem: Agentic memory for llm agents. *arXiv preprint arXiv:2502.12110*.

Hui Yang, Sifu Yue, and Yunzhong He. 2023. Auto-gpt for online decision making: Benchmarks and additional opinions. *arXiv preprint arXiv:2306.02224*.

30

Junwei Yang, Hanwen Xu, Srbuhi Mirzoyan, Tong Chen, Zixuan Liu, Zequn Liu, Wei Ju, Luchen Liu, Zhiping Xiao, Ming Zhang, et al. 2024a. Poisoning medical knowledge using large language models. *Nature Machine Intelligence*, 6(10):1156–1168.

Zhicheng Yang, Jinghui Qin, Jiaqi Chen, Liang Lin, and Xiaodan Liang. 2022. Logicsolver: Towards interpretable math word problem solving with logical prompt-enhanced learning. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 1–13.

Zonglin Yang, Wanhao Liu, Ben Gao, Tong Xie, Yuqiang Li, Wanli Ouyang, Soujanya Poria, Erik Cambria, and Dongzhan Zhou. 2024b. Moose-chem: Large language models for rediscovering unseen chemistry scientific hypotheses. *Preprint*, arXiv:2410.07076.

Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023. React: Synergizing reasoning and acting in language models. In *International Conference on Learning Representations (ICLR)*.

Geyan Ye, Xibao Cai, Houtim Lai, Xing Wang, Junhong Huang, Longyue Wang, Wei Liu, and Xiangxiang Zeng. 2023a. Drugassist: A large language model for molecule optimization. *arXiv preprint arXiv:2401.10334*.

Shaokai Ye, Jessy Lauer, Mu Zhou, Alexander Mathis, and Mackenzie Mathis. 2023b. Amadeusgpt: a natural language interface for interactive animal behavioral analysis. *Advances in neural information processing systems*, 36:6297–6329.

Ruihong Zeng, Jinyuan Fang, Siwei Liu, and Zaiqiao Meng. 2024. On the structural memory of llm agents. *arXiv preprint arXiv:2412.15266*.

Huan Zhang, Yu Song, Ziyu Hou, Santiago Miret, and Bang Liu. 2024a. HoneyComb: A flexible LLM-based agent system for materials science. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 3369–3382, Miami, Florida, USA. Association for Computational Linguistics.

Jiaxin Zhang, Zhongzhi Li, Ming-Liang Zhang, Fei Yin, Cheng-Lin Liu, and Yashar Moshfeghi. 2024b. Geoeval: Benchmark for evaluating llms and multi-modal models on geometry problem-solving. In *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, pages 1258–1276. Association for Computational Linguistics.

Shiqi Zhang, Xinbei Ma, Zouying Cao, Zhuosheng Zhang, and Hai Zhao. 2025. Plan-over-graph: Towards parallelable llm agent schedule. *arXiv preprint arXiv:2502.14563*.

Yu Zhao, Huifeng Yin, Bo Zeng, Hao Wang, Tianqi Shi, Chenyang Lyu, Longyue Wang, Weihua Luo, and Kaifu Zhang. 2024a. Marco-o1: Towards open reasoning models for open-ended solutions. *arXiv preprint arXiv:2411.14405*.

Zihan Zhao, Da Ma, Lu Chen, Liangtai Sun, Zihao Li, Hongshen Xu, Zichen Zhu, Su Zhu, Shuai Fan, Guodong Shen, et al. 2024b. Chemdfm: Dialogue foundation model for chemistry. *arXiv preprint arXiv:2401.14818*.

## A Google's AI Co-scientist: An Illustration of a Scientific Agent System

Below is an illustration of how an LLM-based scientific agent, the Google's AI co-scientist (Gottweis et al., 2025), can function as a multi-agent system to assist researchers in formulating hypotheses, designing experiments, and synthesizing existing literature. The components are structured into three core modules, following our earlier description: (1) Planner, (2) Memory, and (3) Tool Set. This organization highlights how advanced language models can iteratively plan scientific inquiries, maintain and refine long-term reasoning, and leverage external resources. An overview of this system is presented in Figure A.1.

### A.1 Planner

The planner provides the overall reasoning and coordination framework, inspired by the steps of the scientific method:

- **Input parsing and configuration.** The system accepts a user defined research goal (e.g., "Propose novel drug repurposing strategies for acute myeloid leukemia") and parses it into a structured research plan. This process encodes any requirements or constraints in the input, such as accessible experimental assays or specific safety considerations.

- **Specialized agents for task execution.** Several specialized agents operate under the planner's coordination:
  - *Generation agent* drafts preliminary hypotheses or proposals, performing literature exploration and combining prior results with new conjectures.
  - *Reflection agent* reviews each hypothesis, checking for consistency, novelty, correctness, and alignment with known data. It can also simulate potential pitfalls or failure points in each proposal.
  - *Ranking agent* organizes a tournament of hypotheses, making pairwise comparisons to assign an Elo-based quality score. Promising ideas are refined further, while those shown to be contradictory or impractical are filtered out.
  - *Proximity agent* computes a similarity graph for hypotheses, enabling clustering and deduplication. By analyzing semantic and contextual relationships between ideas, the agent groups similar hypotheses and identifies redundant ones, ensuring efficient exploration of the hypothesis space.
  - *Evolution agent* iteratively refines top-rated hypotheses, merging or adapting ideas based on feedback.
  - *Meta-review agent* synthesizes recurring observations such as overlooked evidence or repeated mistakes into meta critique feedback for other agents. It also integrates the top-rated hypotheses and reviews into a coherent, high-level overview for the user.

- **Resource Scheduling.** The *Supervisor agent* manages the entire process, allocating computational resources to each specialized agent based on the complexity of the research goal and the system's progress. Through iterative task dispatch, the planner maintains a sustained, self-improving cycle of reasoning.

### A.2 Memory

Ensuring continuous and coherent multi-step reasoning demands robust mechanisms for storing, retrieving, and updating the system's state. To this end, the system employs a persistent context memory that supports iterative reasoning cycles. This repository houses newly generated hypotheses, expert commentary, external references, and notes from specialized agents, thereby maintaining continuity throughout the computational workflow. When reflection critiques, tournament rankings, or meta-review insights become available, they are appended to the memory, allowing the system to refine its reasoning while preserving the record of past decisions.

Stateful storage further enables long-horizon iterative refinement, where hypotheses evolve incrementally without compromising earlier logical foundations. For instance, partial experimental details or validated findings remain accessible even as new data are integrated, preventing the loss of critical insights. The memory also tracks resource allocation metrics, such as hypothesis generation success rates, to guide the Supervisor agent in dynamically prioritizing tasks. In addition, the system keeps summaries of key results, including top-ranked hypotheses and recurring pitfalls, to streamline knowledge retrieval for both human
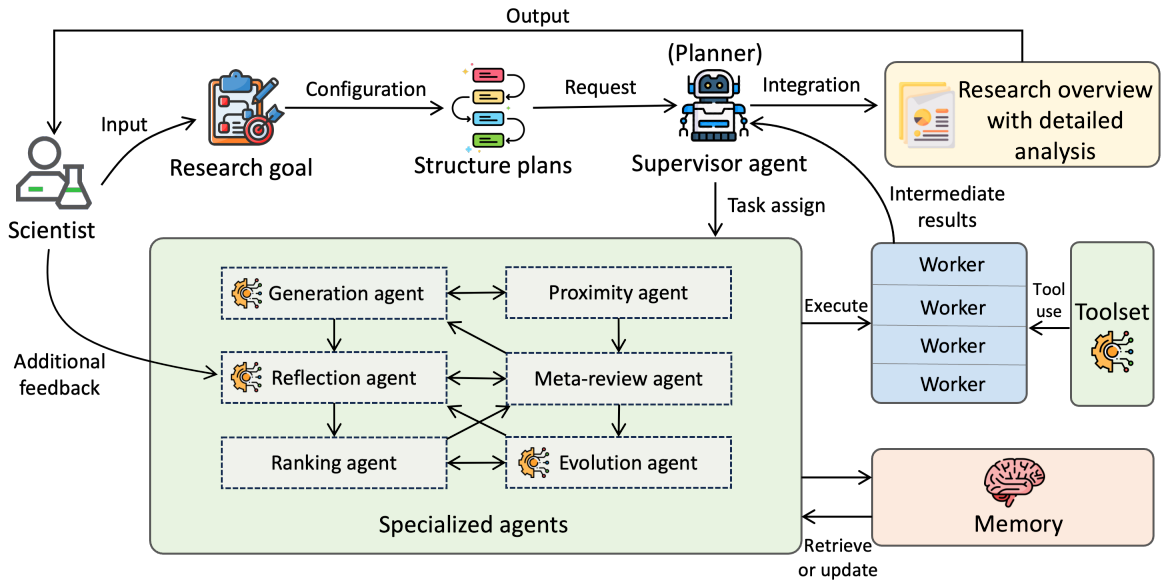
Figure A.1: **The Google's AI co-scientist multi-agent architecture (Gottweis et al., 2025).** The system begins by accepting a natural language research goal from the scientist, which is parsed into a structured research plan. This plan is forwarded to the Supervisor Agent, which evaluates its requirements to allocate computational resources and priority weights to specialized agents. These agents are then organized into a task queue based on their assigned weights. Notably, only the Generation, Reflection, and Evolution agents have access to the toolset. Worker processes execute the queued tasks sequentially, leveraging the expertise of each specialized agent. Finally, the system synthesizes the collected data to generate a comprehensive research summary, including hypotheses and actionable proposals for the user. In the "Specialized Agents" section, grey boxes highlight individual agents, each designed with distinct operational logic and task-specific roles.

users and specialized agents. These mechanisms enable the co-scientist to balance innovation and historical awareness, fostering sustained, cumulative progress toward research goals.

### A.3 Tool Set

Meanwhile, the AI co-scientist expands its capabilities beyond text generation through strategic integration of specialized tools. It leverages search and retrieval systems to query literature databases, online repositories, and user-provided resources, ensuring hypotheses are grounded in existing evidence, avoiding redundancy, and identifying gaps for novel insights.

For domain-specific tasks, the system invokes tailored tools such as AlphaFold (Jumper et al., 2021) for protein structure validation, the Cancer Dependency Map (DepMap) (Institute, 2024) for gene dependency analysis in cancer, and drug libraries for repurposing candidates.

### A.4 Summary of Workflow

As a whole, the multi-agent AI co-scientist operates as follows:

1. Parsing user input into structured plans.

2. Generating and reviewing plausible hypotheses through specialized agents. Notably, users also have the flexibility to refine or expand their requirements by interacting with the agent during the hypothesis-generation process.

3. Conducting an iterative tournament where the top ideas are compared, refined, or strategically combined.

4. Maintaining a context repository for long-horizon memory.

5. Accessing relevant tools and specialized models to verify and refine proposals.

6. Producing a final research overview or set of top-ranked candidates, enabling direct engagement with human researchers for real-world validation.

In this way, the Planner, Memory, and Tool Set modules collectively foster a systematic, self-improving approach to advanced scientific inquiry, leveraging LLMs and related tools to complement and amplify human expertise.

## A.5 Discussion

The multi-agent framework of Google's AI co-scientist offers a powerful approach for automating scientific discovery, particularly across three key problem areas in biomedicine. First, the system has been demonstrated to propose promising drug repurposing candidates for diseases such as acute myeloid leukemia. Second, it has shown potential in discovering novel treatment targets, as illustrated in the identification of epigenetic regulators for liver fibrosis. Third, it has helped to uncover mechanisms of microbial evolution and antimicrobial resistance, recapitulating unpublished findings of novel gene transfer pathways in bacteria.

However, several limitations remain. On the one hand, while the multi-agent design helps isolate errors to specific stages of reasoning, the potential for model hallucination requires careful oversight. Relying on automated reviews, even if tournament-based, does not fully eliminate inaccuracies and overconfident claims. On the other hand, the system's recommendations hinge heavily on the corpus of literature and data it can access. In emerging fields or topics with limited public data, the generated ideas may be too speculative or miss key non-public findings. Additionally, as with any AI-driven method, ethical, legal, and regulatory considerations become paramount when moving from in silico predictions to clinical or large-scale biological testing. Lastly, while the method shows promise in biomedical contexts, its generalizability to other domains remains untested and may require field-specific adaptations.

Looking forward, further enhancements of the co-scientist framework can focus on several complementary directions. First, incorporating increasingly multimodal and domain-specialized AI systems has the potential to accelerate discovery not only in oncology, fibrosis, and antimicrobial resistance, but in an expanding range of biomedical domains. Second, refining methods to detect and mitigate hallucinations—through more transparent agent interactions, robust error-logging, and human verification loops—could make multi-agent pipelines more reliable. Finally, applying similar AI multi-agent architectures to emerging therapies, personalized medicine, and even non-biomedical areas of science may further highlight the versatility of large language models as co-collaborators, potentially reshaping entire research workflows.