

# BioAgents: Democratizing Bioinformatics Analysis with Multi-Agent Systems

Nikita Mehandru, Amanda K. Hall, Olesya Melnichenko,  
Yulia Dubinina, Daniel Tsirulnikov, David Bamman,  
Ahmed Alaa, Scott Saponas, Venkat S. Malladi

## Abstract

Creating end-to-end bioinformatics workflows requires diverse domain expertise, which poses challenges for both junior and senior researchers as it demands a deep understanding of both genomics concepts and computational techniques. While large language models (LLMs) provide some assistance, they often fall short in providing the nuanced guidance needed to execute complex bioinformatics tasks, and require expensive computing resources to achieve high performance. We thus propose a multi-agent system built on small language models, fine-tuned on bioinformatics data, and enhanced with retrieval augmented generation (RAG). Our system, BioAgents, enables local operation and personalization using proprietary data. We observe performance comparable to human experts on conceptual genomics tasks, and suggest next steps to enhance code generation capabilities.

## 1 Main

Large language models (LLMs) have been applied to various domain-specific contexts, including scientific discovery in medicine [45, 49, 56], chemistry [6, 7], and biotechnology [31]. Recent advances in LLMs have spurred their use in bioinformatics [13], a field encompassing data-intensive tasks such as genome sequencing, protein structure prediction, and pathway analysis. One of the most significant applications has been AlphaFold3, which uses transformer architecture with triangular attention to predict a protein’s three-dimensional (3-D) structure from amino acid sequences [2]. Other applications include the use of protein language models in transforming amino acids into embeddings [58].

While LLMs demonstrate impressive capabilities, these models have been found to struggle on complex genomics [4, 29, 39] and bioinformatics code generation [24, 47, 51] tasks with their performance and time to arrive at the correct solution varying significantly with task complexity. For example, ChatGPT solved 97.3% of programming exercises from an introductory bioinformatics course within seven or fewer tries [37]. The model; however, was only able to

solve 75.5%, or 139 out of 184 exercises, on its first attempt. This disparity highlights that while LLMs can assist bioinformatics researchers with introductory data analysis questions, they encounter challenges when tackling more intricate and complex real-world programming, analysis, and research questions, often requiring knowledge of multiple tools, data formats, and analysis techniques.

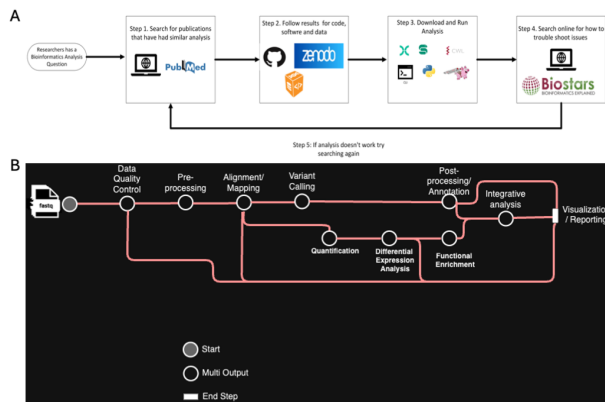


Figure 1: **Supporting Research through Knowledge Graphs and Directed Acyclic Graphs.** A. A knowledge graph showcases the current state-of-the-art for a researcher to start from a research question and navigate through relevant data, tools, and methods to independently run their analysis. B. A typical bioinformatics workflow represented as a Directed Acyclic Graph (DAG), showcasing the intricate dependencies between tasks such as data preprocessing, genome assembly, annotation, and analysis, where each node represents a computational step and edges indicate the flow of data or control.

A frequent challenge for bioinformatics researchers is navigating the complexity of building end-to-end pipelines, which typically requires expertise across multiple domains. Bioinformaticians often mine question-answer platforms like BioStars for similar problems, search for reproducible scientific workflow examples (e.g., Nextflow, WDL or Snakemake) and software containers (e.g., Biocontainers) on GitHub [15, 20, 35], or refer to the methods sections of recently published papers for code. The creation of these workflows require several key steps involving various dependencies, software, compute, storage, data, and vast expertise, including data pre-processing, alignment, and post-processing (as shown in Figure 1). This complexity can present a steep learning curve for newcomers, and poses challenges for bioinformatics experts to stay up-to-date with new techniques [9, 43], as well as with analysis-specific software versions. While established open-source community platforms provide one-off exchanges, they offer limited guidance for researchers trying to develop complex, multi-step

workflows across a network of on-premise and cloud infrastructure [36]. As a result, there is a need for interactive and dynamic tools that can offer continuous support.

To bridge this gap and democratize access to bioinformatics knowledge, we introduce BioAgents – a multi-agent system designed to assist users in designing, developing, and troubleshooting complex bioinformatics pipelines. Recognizing the potential of multi-agent frameworks [26, 44, 46, 53], our system provides an interactive solution that adapts to the ongoing needs of users working in specialized domains [19, 33, 50, 54, 55].

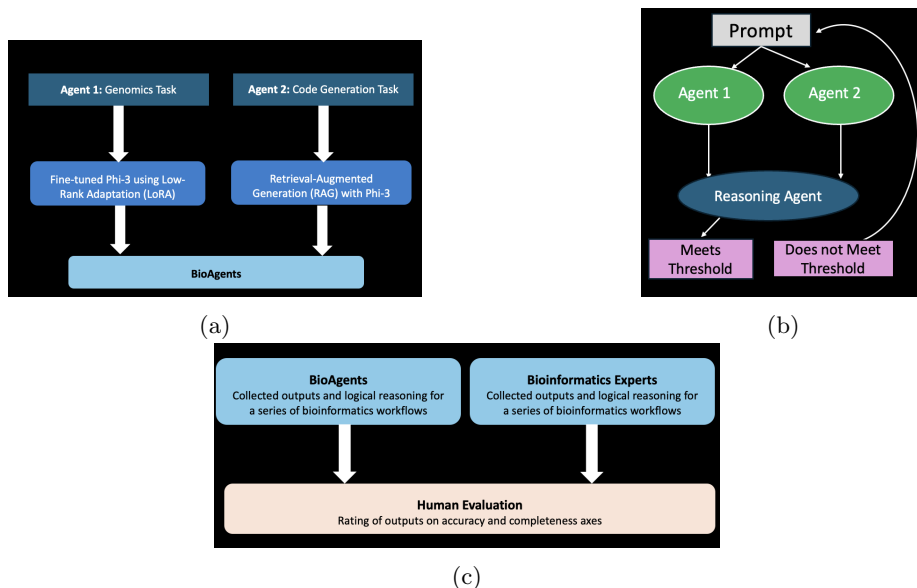


Figure 2: (a) **Two Specialized Agents.** Each specialized agent used Phi-3. The first agent focused on conceptual genomics tasks and was fine-tuned on bioinformatics tools documentation, while the second agent used retrieval-augmented generation (RAG) on workflow documentation. (b) **Overview of BioAgents.** The reasoning agent, a baseline Phi-3 model, processes the outputs from each specialized agent independently and generates the final response. (c) **Comparison of BioAgents' Outputs with Expert Outputs.**

To better understand the challenges faced by practitioners, we analyzed 68,000 question-answer (QA) pairs from Biostars, extracting the associated tags and categorizing each question. The most frequent questions on the platform revolved around tools, specifically bioinformatics software programs and packages, as well as analysis, such as pipeline-related queries focused on RNA-sequencing, alignment, and variant calling. To address the diverse and complex nature of these questions, we employ multiple specialized agents, each tailored to handle specific tasks such as tool selection, workflow generation, and error troubleshooting, enabling a modular and efficient approach to solving bioinformatics chal-

lenges. These insights directly informed the development of our two specialized agents within BioAgents.

While existing multi-agent systems primarily rely on large language models [21, 30], we leveraged a smaller, more efficient language model, Phi-3 [1]. By using a smaller language model, we are able to maintain high performance while significantly reducing computational resources and infrastructure [11, 40]. Avoiding the heavy infrastructure demands associated with larger models, BioAgents is more accessible for local use and efficient real-time applications.

We used the baseline Phi-3 model to build three agents: two specialized agents and Phi-3 as the reasoning agent. Our first agent focused on conceptual genomics tasks, and was fine-tuned on bioinformatics tools documentation from Biocontainers and the software ontology [14, 32]. Our second agent used retrieval-augmented generation (RAG) on nf-core documentation and the EDAM ontology [5, 18, 22]. Figure 2 shows the creation of our two specialized agents, an overview of the BioAgents, and our experimental design, respectively.

## 2 Results

### 2.1 Evaluation across several use cases

We devised three use cases of varying difficulty to evaluate our multi-agent system. These workflows, listed below, were designed to assess both conceptual genomics (analysis steps) and code generation tasks. We recruited bioinformatics experts, and provided them with the same inputs used by the multi-agent system. Each workflow involved completing the conceptual genomics and code generation tasks, providing any additional information needed to aid in answering the user query, and explaining the logical reasoning behind the final output.

#### **Conceptual Genomics and Code Generation Tasks**

##### *Level 1 Tasks (Easy)*

- How would I provide quality metrics on FASTQ files?
- What code or workflow do I need to write to provide quality metrics on FASTQ files?

##### *Level 2 Tasks (Medium)*

- How do I align RNA-seq data against a human reference genome?
- What code or workflow do I need to write to align RNA-seq data against a human reference genome?

##### *Level 3 Tasks (Hard)*

- How can I assemble, annotate, and analyze SARS-CoV-2 genomes from sequencing data to identify and characterize different variants of the virus?

- What code or workflow do I need to write to assemble, annotate, and analyze SARS-CoV-2 genomes from sequencing data to identify and characterize different variants of the virus?

To assess performance, a bioinformatician reviewed both the system and human expert outputs on two axes: 1) accuracy, and 2) completeness. Accuracy was defined as how well the user’s query was answered, while completeness referred to the extent to which the output captured all relevant information in response to the user query. Evaluation results are presented in Figure 3.

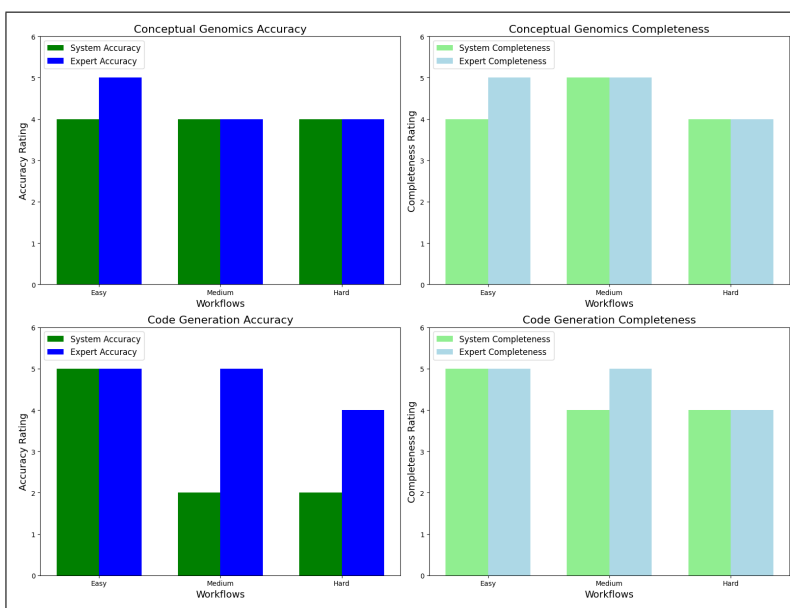


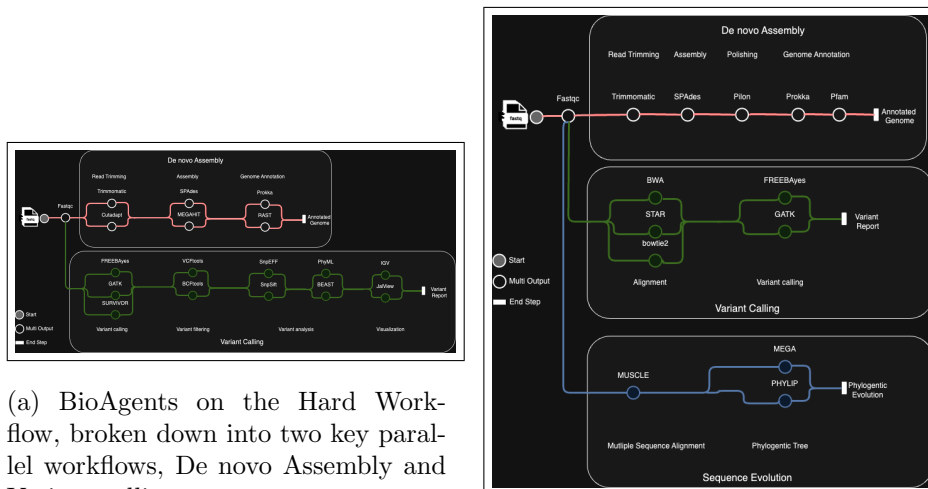
Figure 3: **Comparison of system and expert performance across conceptual genomics and code generation tasks.** The top row evaluates conceptual genomics tasks, with separate panels for accuracy (left) and completeness (right). The bottom row evaluates code generation tasks, similarly split into accuracy (left) and completeness (right). For conceptual genomics tasks, the system demonstrates comparable performance to human experts across all levels of difficulty. In code generation tasks, the system matches expert performance on easier tasks, but shows a decline in accuracy and completeness for medium and hard tasks, highlighting opportunities for improvement in addressing complex challenges.

### 2.1.1 Conceptual Genomics Tasks

BioAgents demonstrated performance on par with experts on conceptual genomics questions across three workflows. This success is largely attributed to our use of Low-Rank Adaptation (LoRA) to fine-tune an agent on the top 50

bioinformatics tools in Biocontainers, including detailed software versions and help documentation. Biocontainers, a widely used bioinformatics service, provides the infrastructure for managing bioinformatics packages and containers, such as conda and docker. Consequently, the system effectively interpreted and responded to these conceptual tasks, achieving human-expert-like performance.

In a challenging workflow question on assembling, annotating, and analyzing SARS-CoV-2 genomes from sequencing data, BioAgents provided a logical series of steps, including obtaining sequencing data, performing quality control, assembling the high-quality reads using a de novo assembler, annotating the assembled genome using tools like Prokka or RAST, identifying and characterizing variants, and constructing a phylogenetic tree [3, 42]. While human experts proposed robust pipelines, they lacked rationales for tool recommendations. BioAgents did occasionally omit steps though, requiring users to fill in the gaps as shown in Figure 4.



(a) BioAgents on the Hard Workflow, broken down into two key parallel workflows, De novo Assembly and Variant calling.

(b) Experts on the Hard Workflow, broken down into three parallel workflows, De novo Assembly, Variant calling, and Sequence evolution.

Figure 4: BioAgents and Experts on the Hard Genomics Workflow

### 2.1.2 Code Generation Tasks

Performance discrepancies emerged in code generation tasks, particularly in workflows of increasing complexity. For easy tasks, BioAgents matched expert accuracy, but sometimes provided false information about tools. For medium tasks, representing end-to-end pipelines like those in nf-core workflows (<https://nf-co.re/pipelines/>), BioAgents struggled to produce complete outputs.

In the most complex workflows, the system failed to generate starter code, instead offering step outlines more similar to a conceptual answer. These limitations were attributed to gaps in the indexed workflows, and a lack of tool and language diversity in the training dataset.

## 2.2 Reliability and Transparency

Two key components are necessary in the deployment of multi-agent systems in highly specialized domains: reliability and transparency. Reliability ensures that the system consistently delivers accurate results, while transparency enables users to understand and trust the system’s decision-making process.

### 2.2.1 Self-Reflection in Agent Systems

Several techniques have been proposed to enable a language model to correct its outputs based on internal evaluation, including: self-consistency [10], self-correction [23, 34], self-evolution [48], self-feedback [27] and self-evaluation [38].

BioAgents incorporated self-evaluation to enhance output reliability, inspired by the idea that agent systems can assess the accuracy of their own outputs [59]. Our reasoning agent assessed the quality of responses against a defined threshold. Outputs scoring below this threshold were reprocessed, with agents independently reanalyzing the prompts before returning results. However, the iterative process revealed diminishing returns, where repeated refinements negatively impacted output quality and might not necessarily lead to improved outcomes.

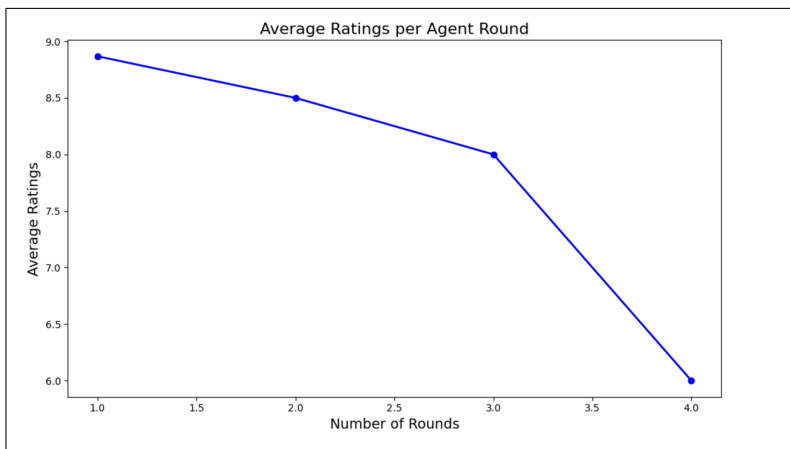


Figure 5: Self-Ratings by Number of Rounds: an inverse correlation between the number of rounds the multi-agent system takes to reach the final answer and the quality of the output’s rating suggests a potential limitation of the iterative processes in multi-agent systems.

### 2.2.2 Collaborative Reasoning and Transparent Guidance

Current applications of LLM agents in domain-specific tasks have struggled in the areas of long-term reasoning, decision-making, and instruction-following [28]. Moreover, their proficiency in addressing practical bioinformatics queries and conducting nuanced knowledge inference remains constrained [12].

In our experimental set-up, both the system and human experts were asked to explain any additional information they would need to better answer users' questions, and the logical reasoning process behind their answers. Our motivation behind this assessment stemmed from various existing reasoning frameworks – chain-of-thought (CoT) [52] and ReAct [57] – used to provide interpretability in LLMs.

In the context of the conceptual genomics medium workflow task, BioAgents explained its rationale for selecting the most suitable RNA-seq alignment tool for mapping against the reference genome, recommending STAR and HISAT2 for their high-throughput and accurate alignments [17, 25]. The multi-agent system described how these alignment tools mapped the RNA-seq reads to the reference genome, enabling the identification of genomic locations within the read. BioAgents also specified the factors influencing its choice of alignment tool, noting the size of the RNA-seq dataset and the user's desired accuracy level as important factors. Generating natural language explanations of model outputs has been shown to improve interpretability, thereby increasing transparency and fostering trust, as users are better able to understand the reasoning behind the model's specific outputs [8, 41].

A key insight from our findings was that our multi-agent system was able to identify additional information that would have improved responses, even in cases where accuracy was lower. For example, in the hard code generation task, BioAgents struggled to generate the necessary workflow code, but was able to identify additional information that could have better answered the user question, specifically more information on the raw sequencing data, the reference genome sequence, software and tool versions, computational resources (e.g., CPU, memory, disk space), and the user's bioinformatics experience. In contrast, although the human experts achieved a higher accuracy score on the same task (four compared to two), they described the limitations of their responses, with one human expert noting that the solutions provided were "cobbled together from searching for tutorials," indicating that it was difficult to identify their information gaps. BioAgents, on the other hand, demonstrated more metacognitive awareness recognizing what it didn't know, and more importantly, additional information it could benefit from [16]. A consistent theme across the human expert responses was an inability to articulate what additional information they would need to improve their answers. This highlights how vast and expansive of a knowledge base is needed to fully answer these questions, which often extends beyond the expertise of one individual.



### 3 Discussion

The reproducibility crisis in computational research highlights the urgent need for systems that can reliably extract, reproduce, and adapt research findings. This challenge is especially pronounced in bioinformatics, where complex workflows often hinder replication and validation efforts. BioAgents, a multi-agent system, offers a promising solution by extracting methods from research papers, generating executable workflows, and integrating human-in-the-loop approaches to improve accuracy and customization.

BioAgents has the potential to facilitate reproducibility by automatically synthesizing workflows from research publications, enabling researchers to replicate experiments, validate results, and adapt analyses to their datasets. By prioritizing transparency and integrating human feedback, BioAgents ensures outputs are both reliable and user-specific. Furthermore, BioAgents can be extended to clinical settings and other scientific domains. In medicine, the system could assist in replicating diagnostic workflows, personalizing treatment recommendations based on patient data, and optimizing clinical trial designs, ultimately enabling more efficient and reliable translational research. In chemistry and physics, BioAgents can automate the replication of experimental protocols and model complex systems, enhancing the reproducibility of results across a wide range of scientific fields.

Collaborative reasoning highlighted key areas for improvement in reasoning, decision-making, and instruction-following. BioAgents effectively identified information gaps, such as tool versions and user experience, which experts often overlooked. Generating natural language explanations of its decisions increased interpretability, fostering user trust. Despite struggles with accuracy in complex tasks, BioAgents demonstrated metacognitive awareness, outlining additional data that could improve results.

The increased reliability, transparency, and trust fostered by BioAgents is particularly valuable for reducing barriers for new bioinformatics researchers. Unlike static question-answer forums, which can provide solutions without clear explanations or insight into the respondent’s reasoning process, BioAgents allows users to not only receive answers, but also provides the underlying supporting information that led to those answers. By sharing the multi-agent system’s reasoning behind its proposed steps, researchers can learn how to replicate those decision-making processes, highlighting the educational value of our approach. Ultimately, the transparency provided by BioAgents not only improves trust in agentic systems but also facilitates knowledge transfer, allowing users to grow and develop their expertise.

In conceptual genomics tasks, BioAgents demonstrated performance comparable to that of human experts, successfully addressing domain-specific challenges. However, areas for improvement were identified in code generation. Specifically:

- **Workflow Scope:** The system’s reliance on nf-core workflows limited diversity. Expanding indexed workflows to include additional sources could

address this gap.

- **Information Retrieval:** Retrieving multiple document matches, rather than relying solely on top-ranked results, could enhance the relevance of generated workflows.
- **Reasoning Agent:** Enhancing the reasoning agent to verify tool versions, ensure executability, and reference source documentation could increase transparency and foster user trust.

One of BioAgents’ key strengths is its ability to support user learning. By linking generated workflows to source documentation and providing support information for each step, BioAgents enables researchers to understand and modify workflows, thereby contributing to their professional development and the broader bioinformatics community. Additionally, when BioAgents requests supplementary information, researchers gain the opportunity to refine results tailored to their specific analyses, while also enhancing their understanding of their own data.

By lowering barriers to compute resources and operating seamlessly in local environments, BioAgents addresses both accessibility and scalability. Its potential extends beyond bioinformatics, offering a model for intelligent systems in other domains facing reproducibility challenges to use the BioAgents framework and train on their own propriety or domain-specific code and documentation.

With targeted enhancements in workflow diversity, retrieval methods, and reasoning capabilities, BioAgents is poised to become a cornerstone in the push for reproducible, transparent, and accessible computational research.

## 4 Methods

### 4.1 Datasets

#### 4.1.1 Biostars

Biostars [35] is an online community platform for the bioinformatics community where users can ask and answer questions related to computational genomics and biological data analysis. We scraped all publicly available data from the site, which included a total of 68,000 question-answer (QA) pairs, up to May 1, 2024. Only answers with at least one user upvote were added to the QA dataset. Tags assigned to each question were extracted, then GPT-3.5 was used to categorize them into one of five categories:

- tool – software programs and packages used for bioinformatics analysis;
- analysis – pipelines and analysis performed in bioinformatics field, such as rna-seq, alignment, variant calling;
- data format – genomics and other -omics data formats;

- programming – programming languages, including wdl, nextflow and snake-make, and operation systems;
- other – for everything else.

#### 4.1.2 Biocontainers

We fine-tuned BioAgents on Biocontainers’ (<https://biocontainers.pro/>) top 50 tools, including versions and documentation. We use the TRS API (<https://api.biocontainers.pro/ga4gh/trs/v2/ui/>) to pull statistics on the top bioinformatics tools, based on download frequency. Then, we grabbed each available docker version of each of those tools. For each docker container, we downloaded the container and outputted command-line help documentation.

#### 4.1.3 Ontologies

We downloaded both the Software and EDAM Ontologies [5, 22, 32] for software and assay description. To convert to JSON-LD we used the JSON or OBO format to extract the name and either the description or definition .

## 4.2 Models

We leveraged a single A100 GPU to perform parameter-efficient fine-tuning (PEFT) on the Phi-3-mini-128-instruct model, optimizing it for bioinformatics tasks. Specifically, we employ the QLoRA technique, which enables fine-tuning with reduced computational overhead by quantizing the model’s layers and training low-rank adapters. This approach is particularly well-suited for large-scale language models like Phi-3-mini, as it retains performance while significantly reducing resource requirements.

Our fine-tuning dataset focuses on the top 50 most commonly used BioContainers tools, along with their associated versions and help documentation, ensuring broad applicability to bioinformatics workflows. Additionally, we added Software Ontology data about the name of the software and its purpose. Training was conducted on Azure Machine Learning, with model configurations limited to a new token count of 1,000 and a temperature of 0.1 to control response diversity and precision.

For our retrieval-augmented generation (RAG) implementation, we integrate OpenAI’s text-embedding-ada-002 for high-quality semantic search. The embeddings are indexed within Azure AI’s search service, optimized to retrieve n-core modules efficiently, and the Sequence Ontology, describing each assay and its purpose. This combination ensures that the system can provide relevant, tool-specific code generation and guidance tailored to bioinformatics workflows.

By combining QLoRA-based fine-tuning and RAG, we achieve a system that balances computational efficiency, domain specificity, and accessibility for researchers in bioinformatics.

## 4.3 Expert Survey

We conducted a survey to elicit expert bioinformatician responses to workflow questions informed by Biostars data, and evaluated this against outputs from BioAgents. Our survey assessed human expert’s reasoning and logic behind their responses. Below we discuss our survey design, respondent recruitment, informed consent process, survey data analysis, and findings.

### 4.3.1 Survey Design

We created a survey using Microsoft Forms to obtain human expert answers and their reasoning behind responses related to translating genomics tasks, and writing subsequent code to analyze data question types. Biostars community forum questions were abstracted and categorized to evaluate common question types, which we found to be around tools and/or analysis. We then created three levels of questions (easy, medium, and hard) increasing in complexity (i.e., number of steps and knowledge required) derived from questions in Biostars.

The survey consisted of 27 questions. We asked eight demographic questions to assess respondent’s education-level, number of years working with bioinformatics tools, data types they worked with regularly, work setting, programming experience, age, gender, and if English was their first language. For the remaining survey questions, respondents were asked to imagine an undergraduate student asked for help on a set of questions, and to provide the logic and steps they would advise the student to take to successfully complete their inquiry. Each of the three question levels had two parts: 1) answering the question asked by the student with corresponding logic/ reasoning, and 2) generating code (in their preferred programming language) with corresponding logic /reasoning. An optional question was asked after each of these questions around what additional information, if any, they would need to answer the student’s question.

### 4.3.2 Recruitment and Informed Consent

We recruited survey respondents through a combination of Microsoft internal community distribution email lists and snowball sampling techniques in August of 2024. The primary inclusion criteria for participation required individuals to have 5+ years of bioinformatics experience with tools/ workflows, intermediate to advanced coding skills, and be at least 18 years of age.

Five bioinformatician experts responded to our survey. The expert respondents ranged in age (25-44 years), gender (2 women, 3 men), and education level (2 masters, 3 doctorates). All respondents reported intermediate to advanced programming experience, 5 or more years of experience working with bioinformatics tools, and English as their first language. Additionally, all respondents reported working with a range of biomedical data types (3 non-human, 3 human, non-clinical, and 4 human, clinical data) and most reported working in an industry setting.

Respondents were excluded if they had higher executive roles in industry to avoid any potential bias. Respondents interested in the study were sent the Microsoft Forms survey to assess their eligibility based on the study criteria, and to review the study informed consent. If they met the eligibility criteria requirements and agreed to the terms outlined in the informed consent, they proceeded to respond to the survey questions. Respondents were compensated \$50 USD for their time via a gift card. This study was reviewed and approval by Microsoft Research Institutional Review Board (ID10950).

### 4.3.3 Qualitative Data Analysis

We conducted a content and thematic analysis of the open-ended questions. For the three question levels (easy, medium, hard) related to the genomics workflows and code generation tasks, we conducted a content analysis of the steps and code across respondents to create a ground-truth dataset to evaluate against the output of BioAgents. For the logic/reasoning and additional information needed questions we conducted a thematic analysis of emerging themes.

### 4.3.4 Expert Findings

We aggregated survey responses from the five expert respondents on Agent 1 and Agent 2 questions across the three levels of difficulty (Figure 2). We created a master list of steps by task, and corresponding code as ground-truth data across all experts' responses to the genomics workflow and code generation questions. We then conducted a thematic analysis for emerging themes based on responses to the logic/reasoning and additional information needed questions.

Expert responses to questions and corresponding logic plus additional information required did not vary across experts on the easy question types. However, for the medium and hard question sets, experts required more logic/reasoning and additional information due to the increased complexity (e.g., number of steps) of the biomedical research tasks necessary to correctly answer the question. This also included the tools and knowledge needed to navigate and implement the correct solution.

## 4.4 Human Evaluation

To assess the performance of the system and experts, we conducted a human evaluation study in which a domain expert scored outputs on two key criteria: **accuracy** and **completeness**. For conceptual tasks, experts assessed whether the system's reasoning and recommendations were consistent with domain knowledge. For code generation tasks, experts verified syntax correctness, tool compatibility, and functionality.

1. **Accuracy (1-5):** The degree to which the output code and steps are correct. The steps and/or software are reasonable and not hallucinated. Where 1 indicates major inaccuracies and 5 indicates full correctness.

2. **Completeness (1-5):** The extent to which the output provides all necessary components or steps required for the task, where 1 indicates significant omissions and 5 indicates comprehensive coverage of code or steps such that the user would be able to implement the answer without searching for additional information.

## 4.5 Supplementary Material

Our benchmarking results indicate that Phi-3 and GPT-4 perform similarly on Biostars QA pairs.

Model	ROUGE-1	ROUGE-2	ROUGE-L	ROUGE-L-SUM
Phi-3.5-mini	0.129	0.015	0.074	0.092
Phi-3.5-MoE	0.129	0.015	0.074	0.091
GPT-4	0.183	0.029	0.103	0.125
GPT-4o	0.122	0.014	0.072	0.091
BioAgents	0.121	0.012	0.071	0.086

Table 1: Benchmarking on 71 Biostars Question-Answer (QA) Pairs

Outputs from BioAgents and our Human Experts for the easy and medium workflows are displayed below:

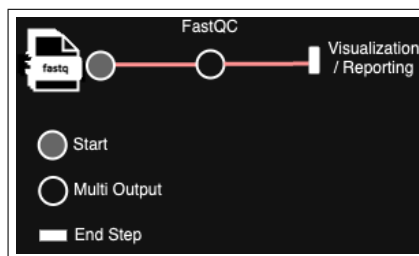


Figure 6: Multi-Agent System on the Easy Workflow

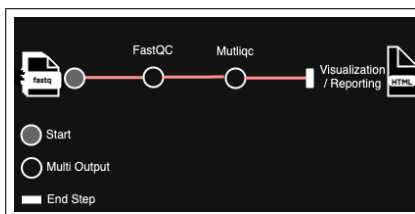


Figure 7: Experts on Easy Workflow

## References

- [1] ABDIN, M., JACOBS, S. A., AWAN, A. A., ANEJA, J., AWADALLAH, A., AWADALLA, H., BACH, N., BAHREE, A., BAKHTIARI, A., BEHL, H., ET AL. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219* (2024).
- [2] ABRAMSON, J., ADLER, J., DUNGER, J., EVANS, R., GREEN, T., PRITZEL, A., RONNEBERGER, O., WILLMORE, L., BALLARD, A. J.,

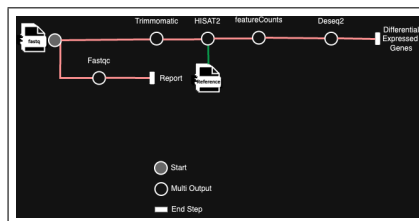


Figure 8: Multi-Agent System on the Medium Workflow

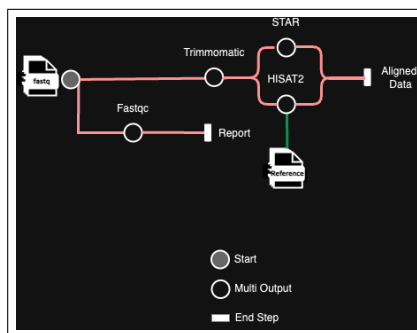


Figure 9: Experts on Medium Workflow

BAMBRICK, J., ET AL. Accurate structure prediction of biomolecular interactions with alphafold 3. *Nature* (2024), 1–3.

- [3] AZIZ, R. K., BARTELS, D., BEST, A. A., DEJONGH, M., DISZ, T., EDWARDS, R. A., FORMSMA, K., GERDES, S., GLASS, E. M., KUBAL, M., ET AL. The rast server: rapid annotations using subsystems technology. *BMC genomics* 9 (2008), 1–15.
- [4] BHARDWAJ, S., AND HASIJA, Y. Chatgpt, a powerful language model and its potential uses in bioinformatics. In *2023 14th International Conference on Computing Communication and Networking Technologies (ICCCNT)* (2023), IEEE, pp. 1–6.
- [5] BLACK, M., LAMOTHE, L., ELDAKROURY, H., ET AL. Edam: the bioscientific data analysis ontology (update 2021), 2022. ISCB Comm J, poster, version 1; not peer reviewed.
- [6] BOIKO, D. A., MACKNIGHT, R., KLINE, B., AND GOMES, G. Autonomous chemical research with large language models. *Nature* 624, 7992 (2023), 570–578.
- [7] BRAN, A. M., COX, S., SCHILTER, O., BALDASSARI, C., WHITE, A. D., AND SCHWALLER, P. Chemcrow: Augmenting large-language models with chemistry tools. *arXiv preprint arXiv:2304.05376* (2023).
- [8] BROWN, N. B. Enhancing trust in llms: Algorithms for comparing and interpreting llms. *arXiv preprint arXiv:2406.01943* (2024).
- [9] CHATTERJEE, A., AHN, A., RODGER, E. J., STOCKWELL, P. A., AND ECCLES, M. R. A guide for designing and analyzing rna-seq data. *Gene expression analysis: methods and protocols* (2018), 35–80.
- [10] CHEN, A., PHANG, J., PARRISH, A., PADMAKUMAR, V., ZHAO, C., BOWMAN, S. R., AND CHO, K. Two failures of self-consistency in the multi-step reasoning of llms. *arXiv preprint arXiv:2305.14279* (2023).

- [11] CHEN, L., AND VAROQUAUX, G. What is the role of small models in the llm era: A survey. *arXiv preprint arXiv:2409.06857* (2024).
- [12] CHEN, Q., AND DENG, C. Bioinfo-bench: A simple benchmark framework for llm bioinformatics skills evaluation. *bioRxiv* (2023), 2023–10.
- [13] CHENG, W., SHEN, J., KHODAK, M., MA, J., AND TALWALKAR, A. L2g: Repurposing language models for genomics tasks. *bioRxiv* (2024), 2024–12.
- [14] DA VEIGA LEPREVOST, F., GRÜNING, B. A., ALVES AFLITOS, S., RÖST, H. L., USZKOREIT, J., BARSNES, H., VAUDEL, M., MORENO, P., GATTO, L., WEBER, J., ET AL. Biocontainers: an open-source and community-driven framework for software standardization. *Bioinformatics* *33*, 16 (2017), 2580–2582.
- [15] DI TOMMASO, P., CHATZOU, M., FLODEN, E. W., BARJA, P. P., PALUMBO, E., AND NOTREDAME, C. Nextflow enables reproducible computational workflows. *Nature biotechnology* *35*, 4 (2017), 316–319.
- [16] DIDOLKAR, A., GOYAL, A., KE, N. R., GUO, S., VALKO, M., LILLICRAP, T., REZENDE, D., BENGIO, Y., MOZER, M., AND ARORA, S. Metacognitive capabilities of llms: An exploration in mathematical problem solving. *arXiv preprint arXiv:2405.12205* (2024).
- [17] DOBIN, A., DAVIS, C. A., SCHLESINGER, F., DRENKOW, J., ZALESKI, C., JHA, S., BATUT, P., CHAISSON, M., AND GINGERAS, T. R. Star: ultrafast universal rna-seq aligner. *Bioinformatics* *29*, 1 (2013), 15–21. Epub 2012 Oct 25.
- [18] EWELS, P. A., PELTZER, A., FILLINGER, S., PATEL, H., ALNEBERG, J., WILM, A., GARCIA, M. U., DI TOMMASO, P., AND NAHNSEN, S. The nf-core framework for community-curated bioinformatics pipelines. *Nature biotechnology* *38*, 3 (2020), 276–278.
- [19] GAO, S., FANG, A., HUANG, Y., GIUNCHIGLIA, V., NOORI, A., SCHWARZ, J. R., EKTEFAIE, Y., KONDIC, J., AND ZITNIK, M. Empowering biomedical discovery with ai agents. *Cell* *187*, 22 (2024), 6125–6151.
- [20] GRUENING, B., SALLOU, O., MORENO, P., DA VEIGA LEPREVOST, F., MÉNAGER, H., SØNDERGAARD, D., RÖST, H., SACHSENBERG, T., O’CONNOR, B., MADEIRA, F., ET AL. Recommendations for the packaging and containerizing of bioinformatics software. *F1000Research* *7* (2019), ISCB-Comm.
- [21] GUO, T., NAN, B., LIANG, Z., GUO, Z., CHAWLA, N., WIEST, O., ZHANG, X., ET AL. What can large language models do in chemistry? a comprehensive benchmark on eight tasks. *Advances in Neural Information Processing Systems* *36* (2023), 59662–59688.



- [22] ISON, J., KALAŠ, M., MÉNAGER, H., WILLIGHAGEN, E., GRÜNING, B., AND IGNARD, A. edamontology/edamontology: Edam 1.25, June 2020.
- [23] KAMOI, R., ZHANG, Y., ZHANG, N., HAN, J., AND ZHANG, R. When can llms actually correct their own mistakes? a critical survey of self-correction of llms. *Transactions of the Association for Computational Linguistics* 12 (2024), 1417–1440.
- [24] KANG, K., YANG, Y., WU, Y., AND LUO, R. Integrating large language models in bioinformatics education for medical students: Opportunities and challenges. *Annals of Biomedical Engineering* (2024), 1–5.
- [25] KIM, D., PAGGI, J. M., PARK, C., BENNETT, C., AND SALZBERG, S. L. Graph-based genome alignment and genotyping with hisat2 and hisat-genotype. *Nature biotechnology* 37, 8 (2019), 907–915.
- [26] LI, G., HAMMOUD, H., ITANI, H., KHIZBULLIN, D., AND GHANEM, B. Camel: Communicative agents for” mind” exploration of large language model society. *Advances in Neural Information Processing Systems* 36 (2023), 51991–52008.
- [27] LIANG, X., SONG, S., ZHENG, Z., WANG, H., YU, Q., LI, X., LI, R.-H., WANG, Y., WANG, Z., XIONG, F., ET AL. Internal consistency and self-feedback in large language models: A survey. *arXiv preprint arXiv:2407.14507* (2024).
- [28] LIU, X., YU, H., ZHANG, H., XU, Y., LEI, X., LAI, H., GU, Y., DING, H., MEN, K., YANG, K., ET AL. Agentbench: Evaluating llms as agents. *arXiv preprint arXiv:2308.03688* (2023).
- [29] LUBIANA, T., LOPES, R., MEDEIROS, P., SILVA, J. C., GONCALVES, A. N. A., MARACAJA-COUTINHO, V., AND NAKAYA, H. I. Ten quick tips for harnessing the power of chatgpt in computational biology. *PLoS Computational Biology* 19, 8 (2023), e1011319.
- [30] M. BRAN, A., COX, S., SCHILTER, O., BALDASSARI, C., WHITE, A. D., AND SCHWALLER, P. Augmenting large language models with chemistry tools. *Nature Machine Intelligence* (2024), 1–11.
- [31] MADANI, A., KRAUSE, B., GREENE, E. R., SUBRAMANIAN, S., MOHR, B. P., HOLTON, J. M., OLMOS, J. L., XIONG, C., SUN, Z. Z., SOCHER, R., ET AL. Large language models generate functional protein sequences across diverse families. *Nature Biotechnology* 41, 8 (2023), 1099–1106.
- [32] MALONE, J., BROWN, A., LISTER, A. L., ISON, J., HULL, D., PARKINSON, H., AND STEVENS, R. The software ontology (swo): a resource for reproducibility in biomedical data analysis, curation and digital preservation. *Journal of biomedical semantics* 5 (2014), 1–13.

- [33] MEHANDRU, N., MIAO, B. Y., ALMARAZ, E. R., SUSHIL, M., BUTTE, A. J., AND ALAA, A. Evaluating large language models as agents in the clinic. *NPJ digital medicine* 7, 1 (2024), 84.
- [34] PAN, L., SAXON, M., XU, W., NATHANI, D., WANG, X., AND WANG, W. Y. Automatically correcting large language models: Surveying the landscape of diverse self-correction strategies. *arXiv preprint arXiv:2308.03188* (2023).
- [35] PARNELL, L. D., LINDENBAUM, P., SHAMEER, K., DALL’OLIO, G. M., SWAN, D. C., JENSEN, L. J., COCKELL, S. J., PEDERSEN, B. S., MANGAN, M. E., MILLER, C. A., ET AL. Biostar: an online question & answer resource for the bioinformatics community. *PLoS computational biology* 7, 10 (2011), e1002216.
- [36] PATEL, K. A beginner’s guide to bioinformatics. *The Biochemist* 45, 2 (2023), 11–15.
- [37] PICCOLO, S. R., DENNY, P., LUXTON-REILLY, A., PAYNE, S., AND RIDGE, P. G. Many bioinformatics programming tasks can be automated with chatgpt. *arXiv preprint arXiv:2303.13528* (2023).
- [38] REN, J., ZHAO, Y., VU, T., LIU, P. J., AND LAKSHMINARAYANAN, B. Self-evaluation improves selective generation in large language models. In *Proceedings on* (2023), PMLR, pp. 49–64.
- [39] SARWAL, V., MUNTEANU, V., SUHODOLSKI, T., CIORBA, D., ESKIN, E., WANG, W., AND MANGUL, S. Biollmbench: A comprehensive benchmarking of large language models in bioinformatics. *bioRxiv* (2023), 2023–12.
- [40] SCHICK, T., DWIVEDI-YU, J., DESSÌ, R., RAILEANU, R., LOMELI, M., HAMBRO, E., ZETTLEMOYER, L., CANCEDDA, N., AND SCIALOM, T. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems* 36 (2024).
- [41] SCHWARTZ, S., YAELI, A., AND SHLOMOV, S. Enhancing trust in llm-based ai automation agents: New considerations and future challenges. *arXiv preprint arXiv:2308.05391* (2023).
- [42] SEEMANN, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* 30, 14 (2014), 2068–2069.
- [43] SHUE, E., LIU, L., LI, B., FENG, Z., LI, X., AND HU, G. Empowering beginners in bioinformatics with chatgpt. *Quantitative Biology* 11, 2 (2023), 105–108.
- [44] SINGH, A., EHTESHAM, A., MAHMUD, S., AND KIM, J.-H. Revolutionizing mental health care through langchain: A journey with a large language model. In *2024 IEEE 14th Annual Computing and Communication Workshop and Conference (CCWC)* (2024), IEEE, pp. 0073–0078.

- [45] SINGHAL, K., AZIZI, S., TU, T., MAHDAVI, S. S., WEI, J., CHUNG, H. W., SCALES, N., TANWANI, A., COLE-LEWIS, H., PFOHL, S., ET AL. Large language models encode clinical knowledge. *Nature* 620, 7972 (2023), 172–180.
- [46] SREEDHAR, K., AND CHILTON, L. Simulating human strategic behavior: Comparing single and multi-agent llms. *arXiv preprint arXiv:2402.08189* (2024).
- [47] TANG, X., QIAN, B., GAO, R., CHEN, J., CHEN, X., AND GERSTEIN, M. B. Biocoder: a benchmark for bioinformatics code generation with large language models. *Bioinformatics* 40, Supplement\_1 (2024), i266–i276.
- [48] TAO, Z., LIN, T.-E., CHEN, X., LI, H., WU, Y., LI, Y., JIN, Z., HUANG, F., TAO, D., AND ZHOU, J. A survey on self-evolution of large language models. *arXiv preprint arXiv:2404.14387* (2024).
- [49] THIRUNAVUKARASU, A. J., TING, D. S. J., ELANGOVAN, K., GUTIERREZ, L., TAN, T. F., AND TING, D. S. W. Large language models in medicine. *Nature medicine* 29, 8 (2023), 1930–1940.
- [50] WANG, L., MA, C., FENG, X., ZHANG, Z., YANG, H., ZHANG, J., CHEN, Z., TANG, J., CHEN, X., LIN, Y., ET AL. A survey on large language model based autonomous agents. *Frontiers of Computer Science* 18, 6 (2024), 186345.
- [51] WANG, Z., ZHOU, Z., SONG, D., HUANG, Y., CHEN, S., MA, L., AND ZHANG, T. Where do large language models fail when generating code? *arXiv preprint arXiv:2406.08731* (2024).
- [52] WEI, J., WANG, X., SCHUURMANS, D., BOSMA, M., XIA, F., CHI, E., LE, Q. V., ZHOU, D., ET AL. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems* 35 (2022), 24824–24837.
- [53] WU, Q., BANSAL, G., ZHANG, J., WU, Y., ZHANG, S., ZHU, E., LI, B., JIANG, L., ZHANG, X., AND WANG, C. Autogen: Enabling next-gen llm applications via multi-agent conversation framework. *arXiv preprint arXiv:2308.08155* (2023).
- [54] XIAO, Y., LIU, J., ZHENG, Y., XIE, X., HAO, J., LI, M., WANG, R., NI, F., LI, Y., LUO, J., ET AL. Cellagent: An llm-driven multi-agent framework for automated single-cell data analysis. *bioRxiv* (2024), 2024–05.
- [55] XIN, Q., KONG, Q., JI, H., SHEN, Y., LIU, Y., SUN, Y., ZHANG, Z., LI, Z., XIA, X., DENG, B., ET AL. Bioinformatics agent (bia): Unleashing the power of large language models to reshape bioinformatics workflow. *bioRxiv* (2024), 2024–05.

- [56] YANG, R., TAN, T. F., LU, W., THIRUNAVUKARASU, A. J., TING, D. S. W., AND LIU, N. Large language models in health care: Development, applications, and challenges. *Health Care Science* 2, 4 (2023), 255–263.
- [57] YAO, S., ZHAO, J., YU, D., DU, N., SHAFRAN, I., NARASIMHAN, K., AND CAO, Y. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629* (2022).
- [58] YIN, H., GU, Z., WANG, F., ABUDUHAIBAIEER, Y., ZHU, Y., TU, X., HUA, X.-S., LUO, X., AND SUN, Y. An evaluation of large language models in bioinformatics research. *arXiv preprint arXiv:2402.13714* (2024).
- [59] ZHUGE, M., ZHAO, C., ASHLEY, D., WANG, W., KHIZBULLIN, D., XIONG, Y., LIU, Z., CHANG, E., KRISHNAMOORTHY, R., TIAN, Y., ET AL. Agent-as-a-judge: Evaluate agents with agents. *arXiv preprint arXiv:2410.10934* (2024).