

# BIOKGBENCH: A KNOWLEDGE GRAPH CHECKING BENCHMARK OF AI AGENT FOR BIOMEDICAL SCIENCE

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Pursuing artificial intelligence for biomedical science, a.k.a. AI Scientist, draws increasing attention, where one common approach is to build a copilot agent driven by Large Language Models (LLMs). However, to evaluate such systems, researchers typically rely on direct Question-Answering (QA) to the LLM itself or through biomedical experiments. How to benchmark biomedical agents precisely from an AI Scientist perspective remains largely unexplored. To this end, we draw inspiration from scientists’ crucial ability to understand the literature and introduce BioKGBench. In contrast to traditional evaluation benchmarks that focus solely on factual QA, where the LLMs are known to have hallucination issues, we first disentangle “Understanding Literature” into two atomic abilities: i) “Understanding” the unstructured text from research papers by performing scientific claim verification, and ii) interacting with structured Knowledge-Graphs for Question-Answering (KGQA) as a form of “Literature” grounding. We then formulate a novel agent task, dubbed KGCheck, using KGQA and domain-based Retrieval-Augmented Generation (RAG) to identify factual errors in existing large-scale knowledge graphs. We collect over two thousand data points for the two atomic tasks and 225 high-quality annotated samples for the agent task. Surprisingly, we find that state-of-the-art general and biomedical agents have either failed or performed inferiorly on our benchmark. We then introduce a simple yet effective baseline, dubbed BKGAgent. On the widely used popular knowledge graph, we discover over 90 factual errors which provide scenarios for agents to make discoveries and demonstrate the effectiveness of our approach.

## 1 INTRODUCTION

Large Language Models (LLMs) are so powerful that they facilitate nearly every aspect of daily life and work right now, even research (Zhao et al., 2023; Baek et al., 2024; He et al., 2023; Zhou et al., 2023). Observing their marvelous successes in text generation (Yu et al., 2022; Celikyilmaz et al., 2020), text summarization (El-Kassas et al., 2021; Gambhir & Gupta, 2017), and other tasks (Jin et al., 2024a; Tang et al., 2023a), along with their consistent failures such as hallucinations (Ji et al., 2023; Yao et al., 2023), one can conclude that LLMs are powerful in certain tasks involving large-scale unstructured data like daily text or images, but relatively powerless when dealing with data-hungry scenarios. As such, researchers then construct AI agents (Wu et al., 2023b; Tian et al., 2023) assisting LLMs with external tools to extend the capabilities of LLMs. These attempts are fruitful in many fields, including autonomous computers (Steiner, 2008), shopping web-agent (Lee & Liu, 2004), code development (Dalle & David, 2004), society simulation (Drogoul & Ferber, 2018; Lan et al., 2023), etc. A natural subsequent attempt is to develop AI agents to simulate scientists, aiding or even taking over the process of scientific discovery (Baek et al., 2024).

As in Figure 1, existing attempts can be grouped into two categories: i) to build an AI agent for a specific task, such as Question Answering (QA) in a specific domain (Zhang et al., 2018); ii) to encompass multiple AI agents to formulate a multi-agent system as the copilot of scientists, automating certain scientific activities, such as experiment result analysis (Bi et al., 2023; Wang et al., 2023b).

Literature review is the most critical ability that a scientist should possess (Snyder, 2019; Thomas et al., 2020). It does not only involve reading and memorizing, but also requires scientists to

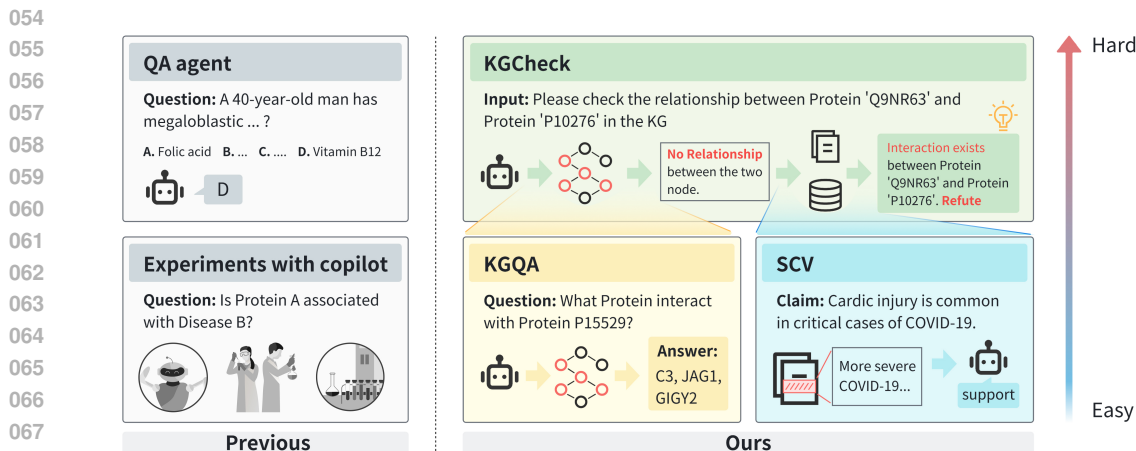


Figure 1: **(Left)** Previous benchmarks for domain-specific AI Agents either focus on the low-level tasks like question answering or are embedded in a complicated pipeline as a scientist copilot. **(Right)** We close the gap by constructing a knowledge graph checking task that consists of two atomic sub-tasks: Knowledge Graph Question Answering (KGQA) and Scientific Claim Verification (SCV), to provide a better evaluation of AI Agents in biomedical science domain.

understand and critically analyze. Researchers and scientists widely spend a significant amount of time in reading recent works. To save human efforts in scientific discovery, it is necessary for AI scientists to be able to accurately understand and analyze the existing research. Many researchers have dedicated to literature understanding in AI agents (Cai et al., 2024; Li et al., 2024), while a systematic evaluation system is missing and even underexplored. The current finest evaluation system (Cai et al., 2024) consists of multiple-choice questions extracted from literature, which cannot fully reveal the underlying reasoning regime of an agent’s success or failure, leaving no clue for future advancement nor indicating whether the agent understands the reasoning rationale or merely memorizes data patterns.

On the other hand, another crucial research direction is to help AI agents capture the underlying logic of literature through domain-specific Knowledge Graphs (KGs) (Abu-Salih, 2021; Kejriwal, 2019). KGs store massive knowledge triples in a graph-structured format (Hogan et al., 2021; Alqaaidi & Kochut, 2024), complementing LLMs with external knowledge while providing frameworks for interpretation and reasoning (Meyer et al., 2023). However, manually constructing such KGs is both intellectually and physically intensive. These domain-specific KGs require annotators with profound domain-specific knowledge, leading to high costs to create or maintain the knowledge graphs. As such, we observe that the existing and well-known biomedical KGs (Santos et al., 2022; Chandak et al., 2023) are not fully reliable due to outdated information. We attribute such discrepancy to the static nature of KGs, which lack mechanisms for dynamic updates to align with the evolution of external knowledge sources.

In this paper, we propose a novel agent evaluation benchmark BioKGBench to address both challenges simultaneously. As in Figure 1 (right), the ultimate goal of our benchmark is to verify the correctness of nodes and triples in the knowledge graph based on various information, including papers and well-maintained databases. We dub this task Knowledge Graph Checking (KGCheck). Agents need to first query the information recorded on the KGs as directed, then cross-reference this information with external literature or databases to combat hallucinations. This task evaluates the agents’ capacities to both process and understand structured data (like KGs) and unstructured data (like literature). It is worth mentioning that the process of verifying knowledge within KGs closely mirrors the methodology of human scientific research, including database queries and extensive literature reviews. This similarity not only underscores the task’s relevance to real-world scientific inquiry but also provides intriguing insights. Furthermore, we decompose this task into two more atomic subtasks: Knowledge Graph Question Answering (KGQA) and Scientific Claim Verification (SCV), enable a more detailed evaluation of the agents’ capabilities in processing and understanding of structured and unstructured data, respectively.

We extensively analyze existing AI agents on our benchmark and find that none of the existing agents can accomplish our tasks without moderate adaptation. Therefore, we introduce our agent BKGAgent, the first agent framework to interact with external knowledge graphs as well as research papers. Experiments demonstrate fascinating results that our agent is capable of discovering real conflicts in the existing large-scale datasets. Within 225 professional-annotated data in the Clinical Knowledge Graph (CKG) (Santos et al., 2022), our agent BKGAgent successfully identified some conflicting or missing pairs. This evidence further supports the academic value of our agent by providing researchers with a tool to update their own knowledge bases, offering substantial potential in both academic and commercial markets.

## 2 RELATED WORK

**Science Agent.** The swift progression of large language models (LLMs) has catalyzed the widespread deployment of intelligent agents across diverse fields, notably within the science domain. Notable examples include ChemCrow (Bran et al., 2023) and Coscientist (Boiko et al., 2023b) in the field of chemistry, DoInstruct (Bi et al., 2023) in ocean science, and GeneGPT (Jin et al., 2024b), Almanac (Zakka et al., 2024), MedAgents (Tang et al., 2023b) in the biomedical domain, etc. Among them, biomedical agents, in particular, have garnered significant attention due to their critical importance. Biomedical agents (Gao et al., 2024) impact areas ranging from hybrid cell simulation (Xiao et al., 2024), the design of cellular circuits (Chandrasekaran et al., 2024) to the development of new therapies (Zhenzhu et al., 2024) and so on. We posit that biomedical agents will emerge as a focal point of research. However, the current benchmark in this field remains inadequate. For instance, MedAgents is evaluated in MedQA (Zhang et al., 2018), MedMCQA (Pal et al., 2022), PubMedQA (Jin et al., 2019), relying heavily on inherent knowledge of LLMs, which leads to hallucinations easily. Our proposed BioKGBench is a dynamic benchmark that evaluates the capabilities of agents in utilizing external tools and knowledge retrieval, thereby addressing this gap.

**Agent Benchmark.** As agents are progressively applied across various domains, the urgency to construct corresponding benchmarks is escalating. Currently, the majority of benchmarks for evaluating agents adopt the approach of evaluating LLM-as-Agent (Liu et al., 2023c), linking LLMs to external frameworks to assess their performance on specific tasks. For instance, AgentBench (Liu et al., 2023c) is a general benchmark for evaluating an agent’s reasoning and decision-making capabilities, SWE-bench (Jimenez et al., 2023) assesses an agent’s proficiency in software engineering, and AgentClinic (Schmidgall et al., 2024) examines an agent’s performance in a simulated clinical environment. However, a benchmark in AI Scientist perspective remains largely unexplored. Our benchmark originates from this perspective, taking the processing and understanding of large-scale data scenarios as the entry point, representing an initial attempt in this direction.

**Agent Integrating LLMs and KGs.** The collaborative use of LLM and KG has become one of the leading methodologies in contemporary agent design, aimed at alleviating uncertainties stemming from the intrinsic mechanisms of LLMs (Pan et al., 2024; Chen et al., 2023a; Yang et al., 2023c). This paradigm not only capitalizes on the generalization ability of LLMs but also employs KGs as an external, trustworthy, and structured data source, thereby achieving reasoning proficiency that strikingly emulates human intellect (Pan et al., 2024). For instance, StructGPT (Jiang et al., 2023) boosts an LLM’s performance on general questions by tapping into the information from a supplied KG. Similarly, KG-Agent (Jiang et al., 2024b) leverages knowledge from KGs, synthesizing instruction data for fine-tuning an open-sourced LLM, thereby achieving competitive performance on general question-answering tasks. However, to our knowledge, while this paradigm has been widely applied to the general question-answering area, its potential remains untapped in the biomedical field. BKGAgent, hence, is poised to fill this gap.

## 3 BIOKGBENCH

Here, we present our benchmark in detail. As aforementioned, one key ability of “AI Scientists” is to understand domain knowledge. However, current LLM-driven agent systems inevitably suffer from hallucinations as a consequence of the statistical nature of LLMs along with the lack of scientific training data compared to data from daily scenarios. We notice that a recent trend in research is to use AI agents to leverage external tools to address these limitations (Bran et al., 2023; Bi et al., 2023).

Table 2: Comparison with existing well-known benchmarks.

Benchmark	Domain	Dataset Composition	Multi-Turn	Environmental Interaction
MMLU (Hendrycks et al., 2020)	57 subjects	QA	✗	✗
MATH (Hendrycks et al., 2021)	math	QA (including solution)	✗	✗
PubMedQA (Jin et al., 2019)	biomedical science	QA	✗	✗
SWE-bench (Jimenez et al., 2023)	software engineering	Issue text, codebase, gold patch, tests	✗	✓
MT-Bench (Zheng et al., 2023)	writing, math, knowledge	votes, conversations	✓	✓
AgentBench (Liu et al., 2023c)	LLM-as-Agent	8 real-world tasks	✓	✓
<b>BioKGBench (ours)</b>	biomedical science, LLM-as-Agent	QA, KGQA, KG, literature	✓	✓

Drawing inspiration from this, we design two atomic abilities to evaluate AI scientists, i) Knowledge Graph Question Answering (KGQA) aiming to address the hallucination issue by grounding the knowledge with structured knowledge graphs; and ii) Scientific Claim Verification (SCV) based on retrieved text from peer-reviewed research papers. In addition, we propose an encompassing task combining these two atomic abilities, to perform Knowledge Graph Checking (KGCheck) as shown in Figure 1. The motivation behind this stems from our interviews with experts from biomedical domains. Their answers to the question “What is the most expected AI agent you would like to use in your daily research?” often included an AI agent that helps in extensive literature review and claim verification. We report the statistics over the scopes of knowledge search, including knowledge graphs and academic literature, in Table 1 (Cf. Appendix A for more details). As shown in Table 2, compared to existing well-known benchmarks, BioKGBench features:

- **setting**: evaluating LLMs as agents through multi-round interactions with the environment to assess their ability to process and understand large-scale biomedical data.
- **data**: a diverse dataset of structured and unstructured data, allowing agents to derive knowledge from heterogeneous sources and make discoveries.

### 3.1 ATOMIC ABILITY

#### 3.1.1 KNOWLEDGE GRAPH QUESTION ANSWERING

This atomic task in the benchmark is to evaluate the agents’ ability to interact with structured Knowledge Graph Question Answering as a grounding of academic literature. Without loss of generality, we choose Clinical Knowledge Graph (CKG) (Santos et al., 2022) as the source of our data, which is one of the most popular large-scale knowledge graph databases in the biomedical domain. CKG is a knowledge graph database with data imported from diverse biomedical databases, aimed at streamlining automated knowledge discovery through the graph’s extensive information.

As the original database is unnecessarily large, we focus on a sub-graph to mitigate the challenge while preserving all relevant information.

Starting from the origin of CKG—protein, we select the sub-graph to contain exactly 12 categories of biological entities, as indicated in Figure 2. Thus, the sub-graph consists of 484,955 entities (nodes) across 12 categories (Biologically defined) and 18,959,943 relationships (edges) of 18 types, with each type consisting of relationships between a unique pair of entity categories (Cf. Appendix A.2 for more details).

After the sub-graph is ready, we construct the question set for the Question Answering (QA) database in two steps. We first handcraft question templates by selecting biomedical fields and pinpointing

Table 1: Statistics of our BioKGBench.

Task	Main Metrics	Scope	Data		
			Dev	Test	All
KGQA	F1	KG	60	638	698
SCV	Acc.	Text (T)	120	1,265	1,385
KGCheck	EM	KG + T	20	205	225

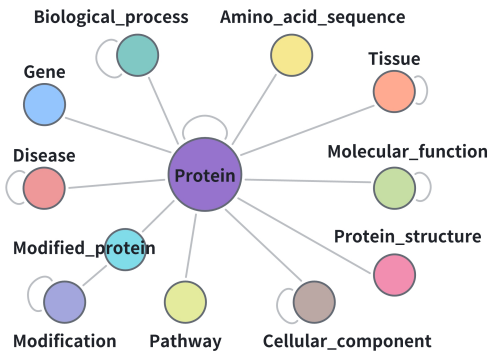
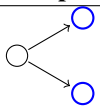
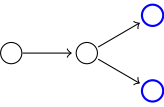
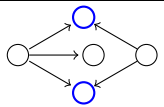


Figure 2: The sub-graph of the Clinical Knowledge Graph (CKG) retains 12 types of nodes and 18 kinds of relationships.



Table 3: Statics of three different reasoning types of KGQA dataset.

Reasoning Type	Graph	Example Question	Question Types	%
One-hop		What proteins does the protein O94842 act on?	8	56.0
Multi-hop		What diseases are associated with the protein encoded by the gene KCNS1?	4	28.7
Conjunction		Which pathway are the proteins P02778 and P25106 both annotated in?	4	15.3

entities and relations in the CKG. Natural language questions were constructed in various formats, ensuring their accuracy through peer reviews and expert consultations. We then expanded our dataset with autogenerated questions by matching CKG data to constructed QA templates, resulting in the generation of 698 questions across three reasoning types and 16 question categories (refer to Table 3).

In this task, we outfit LLMs with a set of atomic KG-querying tools and ask them to answer biomedical questions by querying the provided KG. The responses will be compared with the gold answers and evaluated using the F1 score, where the gold answer to the input question is typically characterized by a set of KG entities. It is noteworthy that our KGQA is built upon a biomedical KG rather than a common sense KG, [with the two adopting different data models. This difference is one of the reasons why KBQA methods cannot be directly applied \(Cf. Appendix C.3 for more details\).](#) This task enables the development of assessing the robustness and tool learning ability of agents built upon various LLMs, and hopefully, it would aid in guiding agents to leverage the extensive biomedical knowledge within the KG, thereby propelling scientific discovery.

### 3.1.2 SCIENTIFIC CLAIM VERIFICATION

This task is designed to evaluate LLMs’ understanding of unstructured text from research papers in a retrieval-augmented generation manner. Following the definition in (Wadden et al., 2020), the task is to identify evidence related to the claim from the research literature and give a verdict of “Support”, “Refute”, or “NEI” (Not Enough Information) based on it. We reconstruct two high-quality biomedical datasets, PubMedQA (Jin et al., 2019) and SciFact (Wadden et al., 2020), into one dataset for SCV, yielding a corpus constituted of abstracts derived from 5,664 scholarly articles, alongside a dataset comprising 1,385 biomedical claims, as shown in Table 4.

Table 4: Examples of reconstructed dataset for SCV, where data from PubMedQA is converted from QA to declarative claims. “NEI” stands for “Not Enough Information”.

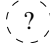
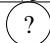
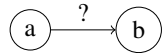
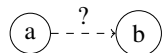
Example Claim	Label	%
A deficiency of folate increases blood levels of homocysteine.	Support	65.2
Therapeutic anticoagulation in the trauma patient is safe.	Refute	33.1
Sternal fracture in growing children is a rare and often overlooked fracture.	NEI	1.7

## 3.2 AGENT TASK

Building upon the atomic abilities, we propose a novel and comprehensive task, KGCheck. This task necessitates the initial application of the tool-query functionality to extract information from the KG. Subsequently, it employs the RAG approach or database access to procure evidence pertaining to the queried information, facilitating a determination of either “Support” or “Refute”. This methodology enables agents to scrutinize the knowledge encapsulated within a large-scale KG, a venture of particular importance considering the prevalence of inaccuracies within numerous datasets, including prominent ones such as ImageNet (Deng et al., 2009).

For this task, we collect 225 high-quality annotated data, as illustrated in Table 5. Given the massive data encapsulated within KGs via triples, we delineate the inspection process into two distinct

Table 5: Four different checking types of KGCheck.

Check Type	Graph	%	Support: Refute
Node	Existence 	20.0	71.0:29.0
	Attribute 	24.4	
Triple	Existing 	25.8	46.4:53.6
	Potential 	29.8	

categories: single-node and triple-based. The single-node inspection is divided into node existence and attribute value assessments, while the triple inspection encompasses scenarios with and without edges between two nodes:

- **Existence:** We note that databases may excise entries during updates due to inaccuracies or redundancies, whereas KGs remain static post-construction, similar to LLMs in some respects. If nodes corresponding to obsolete entities persist in the KG, the label is “Refute”; if they are congruent with real-time updated external databases, the label is “Support”.
- **Attribute:** Our KG is characterized by high information density, with each node and edge encapsulating numerous attribute values, which we scrutinize for accuracy and completeness.
- **Existing Relationship:** We check whether existing edges contradict information from external, real-time updated databases and literature. If external knowledge corroborates the relationship, the label is “Support”; conversely, it is “Refute”.
- **Potential Relationship:** If a relationship is confirmed by databases or literature but is not represented in the KG, the label is “Refute”; otherwise, it is “Support”.

Despite utilizing the latest databases (as of May 2024), we identified errors within the KG, evidenced by 96 “Refute” annotations. These data are valuable and provide scenarios for agents to comprehend knowledge from heterogeneous sources and make **discoveries**.

### 3.3 BKGAGENT: A SIMPLE BASELINE

We propose a **biomedical knowledge-graph agent** (BKGAgent), as shown in Figure 3. It’s a multi-agent framework based on *lang-graph* (Chase, 2023), capable of retrieving information from knowledge graph and cross-validating its correctness with multiple information sources. Our framework is comprised of three agents: the team leader for the progress control, the KG agent for information retrieval from KG, and the validation agent for checking the correctness of the information from KG. This setup simulates the workflow of a human research team, where a leader supervises the assistants’ work and makes the final decision based on their feedback. Additionally, the tool executor is solely responsible for executing functions, and is not based on LLMs.

When a user assigns a task, the leader initially breaks down the task and announces the plan. Then the KG agent is activated to retrieve task-related information from the KG. This involves specifying the tool and its arguments to the tool executor, interpreting the tool result, and communicating it back to the leader. After that, the validation agent is called to verify the information with a workflow similar to that of the KG agent. Finally, the leader will draw a conclusion and return it to the user.

<sup>1</sup><https://gptr.dev/>

Table 6: Comparison of capabilities for BKGAgent and other frameworks.

Framework	MA	KGq	IR
HuggingGPT (Shen et al., 2023)	✓	✗	✗
OpenAgents (Xie et al., 2023)	✓	✗	✗
AgentVerse (Chen et al., 2023b)	✓	✗	✗
Xagent (Team, 2023)	✓	✗	✗
BabyAGI (Yoheinakajima)	✓	✗	✗
MedAgents (Tang et al., 2023c)	✓	✗	✗
gpt-researcher <sup>1</sup>	✓	✗	✗
BDAGENT (Roohani et al., 2024)	✗	✗	✓
<b>BKGAgent(ours)</b>	✓	✓	✓

BDAGENT=BioDiscoveryAgent; MA=multi-agent; KGq=KG-query; IR=information retrieval

324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377

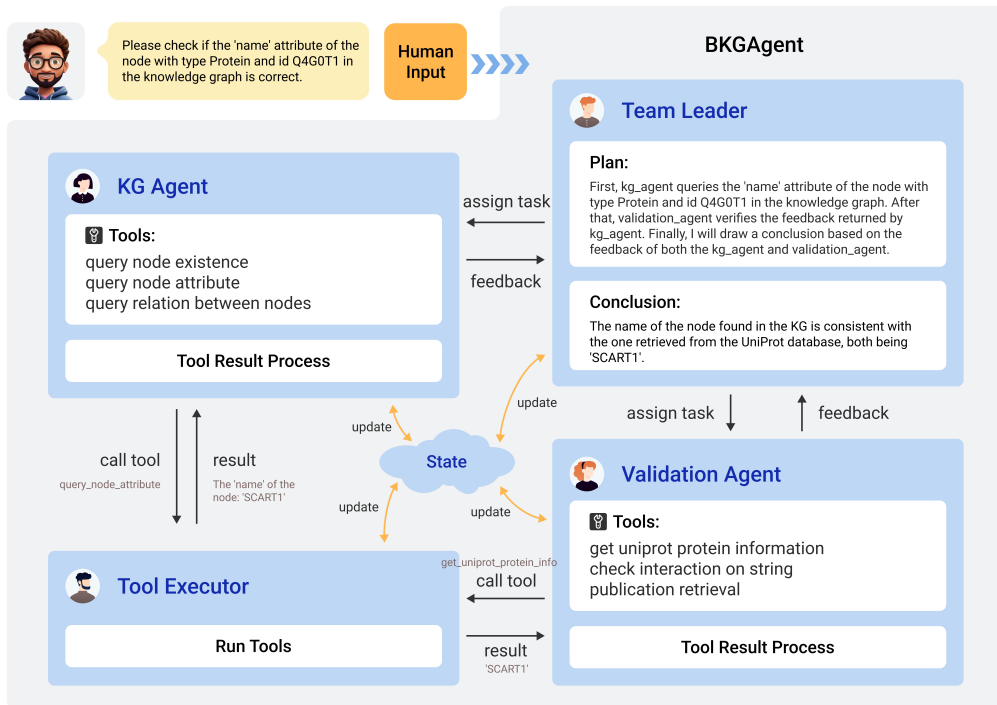


Figure 3: Framework of our BKGAgent.

BKGAgent possesses fundamental capabilities for grounding heterogeneous biomedical knowledge, including knowledge graph queries, database queries, and retrieval-augmented generation (RAG) of literature. In contrast, as illustrated in Table 6, many other frameworks struggle to achieve comparable effectiveness in biomedical information retrieval and verification due to their limited capacity to access knowledge graphs and biomedical data, as well as the unreliability of information sourced from the web.

## 4 EXPERIMENTS

### 4.1 MAIN RESULTS AND ANALYSIS: ATOMIC ABILITIES

**Metrics.** For KGQA, we adopt three metrics: F1, Exact Match (EM), and Executability. For SCV, we adopt three metrics: Accuracy, Right Quotes, and Error Rate. Specifically, “Executability” refers to the success rate of the agent providing an answer within 15 turns, “Right Quotes” indicates the success rate of retrieving matching text from the whole corpus through RAG, and “Error Rate” refers to the frequency with which the agent fails to make a verification.

The detailed experimental results of atomic abilities evaluation on LLMs are shown in Table 7, and we summarize our key findings as follows:

- **Disparity between open-source and commercial API models.** Commercial API models like GPT-4 and GLM-4 generally outperform open-source models in several key metrics. GPT-4, for example, consistently achieves higher scores in both KGQA and SCV tasks, highlighting the advantage of proprietary training techniques and larger computational resources.
- **Strong performance of open-source large models.** Some large OSS models, such as Llama-3-70B-Instruct, perform competitively, sometimes surpassing API models in specific metrics. Llama-3-70B-Instruct, in particular, excels in KGQA executability, suggesting that optimized training can enable open-source models to rival or exceed commercial counterparts.
- **Model parameters do not always correlate with better performance.** In the OSS (Medium) and OSS (Small) categories, smaller models like Llama-3-8B-Instruct sometimes outperform larger models like Qwen1.5-32B-Chat in SCV tasks, indicating that model architecture, training data

Table 7: Test set (standard) results of two easy tasks: KGQA, SCV. **Bold/underline** and **red/blue** indicate the best and second in the subgroup and overall.

LLM Type	Models	KGQA			SCV		
		F1	EM	Executability	Accuracy	Right Quotes	Error
API	GPT-4 (OpenAI, 2023a)	<b>81.8</b>	<b>79.2</b>	<b>88.4</b>	83.9	<b>87.7</b>	<b>0.4</b>
	GLM-4 (Du et al., 2022)	72.4	70.4	82.7	<b>86.9</b>	86.5	0.6
OSS (Large)	Qwen1.5-72B-Chat (Bai et al., 2023)	74.7	72.2	<u>96.1</u>	<u>85.7</u>	<u>83.3</u>	<b>0.1</b>
	Llama-3-70B-Instruct (AI@Meta, 2024)	<b>80.7</b>	<b>77.8</b>	<b>97.0</b>	<b>85.9</b>	<b>86.6</b>	<u>0.2</u>
	DeepSeek-LLM-67B-Chat (Bi et al., 2024)	69.6	66.8	86.3	76.6	82.6	0.4
OSS (Medium)	Qwen1.5-32B-Chat (Bai et al., 2023)	<u>64.6</u>	<b>62.1</b>	<b>83.0</b>	<b>79.7</b>	<b>83.0</b>	<u>0.4</u>
	Qwen1.5-14B-Chat (Bai et al., 2023)	<b>66.0</b>	61.6	78.7	66.1	67.4	<b>0.2</b>
	Baichuan2-13B-Chat (Yang et al., 2023a)	43.7	42.0	<u>82.2</u>	26.3	35.8	<b>0.5</b>
OSS (Small)	Llama-3-8B-Instruct (AI@Meta, 2024)	<b>54.7</b>	<b>51.3</b>	<b>84.8</b>	<b>78.5</b>	<b>83.3</b>	<b>0.5</b>
	Qwen1.5-7B-chat (Bai et al., 2023)	44.5	40.3	77.9	72.5	39.1	2.2
OSS (MoE)	Mixtral-8x7B-Instruct-v0.1 (Jiang et al., 2024a)	<b>70.1</b>	<b>67.9</b>	<b>84.7</b>	<b>77.8</b>	<b>82.5</b>	<u>2.3</u>
	Starling-LM-alpha-8x7B-MoE-GPTQ (Zhu et al., 2023)	12.4	10.9	30.7	<u>55.0</u>	56.2	<b>0.1</b>
	Qwen1.5-MoE-A2.7B-Chat (Bai et al., 2023)	<u>28.7</u>	<u>26.7</u>	<u>71.9</u>	<u>55.0</u>	<u>57.8</u>	3.0

quality, and fine-tuning strategies significantly impact performance. Notably, Qwen1.5-14B-Chat outperforms Qwen1.5-32B-Chat in KGQA, suggesting the latter’s pre-training may be insufficient.

- **Domain-specific models lack transferability.** DeepSeek-LLM-67B-Chat excels in mathematical problems (Bi et al., 2024), but underperforms in biomedical-related tasks, highlighting its lack of cross-domain transferability. This suggests that specialization in one area may compromise generalizability.
- **Inconsistent performance of MoE models.** While Mixtral-8x7B-Instruct-v0.1 performs well in both KGQA and SCV tasks, other MoE models like Starling-LM-alpha-8x7B-MoE-GPTQ and Qwen1.5-MoE-A2.7B-Chat show significantly lower scores. This inconsistency suggests that the effectiveness of MoE models heavily depends on the implementation and integration of the expert models. Additionally, Mixtral-8x7B-Instruct-v0.1, though strong in main metrics, struggles with controlling response format, indicating that individual expert models still require improvement.
- **Biomedical knowledge embedded in model parameters.** The new metric “Right Quotes” for SCV assesses the alignment of retrieved quotes with ground truth evidence. Some models, such as GLM-4, Qwen1.5-72B-Chat, and Qwen1.5-7B-Chat, exhibit higher accuracy metrics than “Right Quotes” metrics. This suggests these models can accurately assess input claims even without sufficient literature evidence, indicating they possess specialized biomedical knowledge.

**Further Analysis.** We also conduct an ablation experiment on three scopes of RAG, as shown in Figure 4, where ‘all’ refers to the abstract of 5,664 articles, ‘partial’ denotes the 1,888 abstracts containing ground truth evidence of claims, and ‘match’ corresponds to the abstracts of the ground truth evidence for the claims. Interestingly, we observe an unexpected phenomenon where the model’s performance in the ‘match’ setting only increases in terms of the right quotes metric, while the accuracy metric decreases. In the ‘all’ setting, we initially anticipated interference from irrelevant literature, but the accuracy metric instead improved. This suggests that there is a potential connection among the extensive literature, where large models exhibit a form of “**analogical reasoning**”. This provides us with insights for conducting extensive literature research in simulating human scientific research.

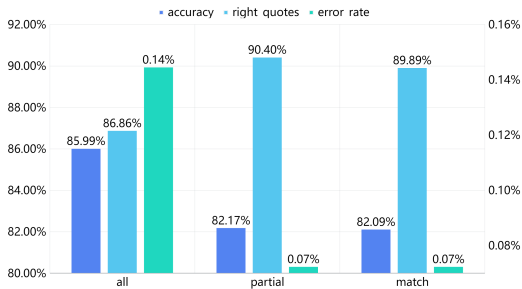


Figure 4: Llama-3-70B-Instruct’s performance in RAG across different scopes of literature.

## 4.2 MAIN RESULTS AND ANALYSIS: BKGAGENT

**Evaluation Setup.** As mentioned in 3.3, most agent frameworks fail in KGCheck, which further highlights that KGCheck is a novel and challenging task. It requires agents to first query the knowledge graph, followed by verification through database searches or RAG of literature. Consequently, agents lacking capabilities for KG querying or information retrieval verification cannot complete this task. Therefore, for agents capable of querying knowledge graphs, we selected KG-Agent from (Liu et al., 2023c) as representative; for general agent frameworks, we chose the most prominent ones,

Table 8: The performance of different agents built on GPT-4o (Hel) in executing the KGCheck task. **Bold/ underline** indicate the best and second, respectively. All scores are on a percentage scale.

Agent	Process					Result	
	Understanding	Reasoning	Efficiency	KG process	Information retrieval	Average	EM
<b>Baselines</b>							
AgentBench-KG agent	81.0	81.0	81.0	77.1	97.1	83.4	56.6
AutoGPT	<b>99.0</b>	78.0	84.4	96.6	59.5	83.5	39.5
AutoGen	95.6	45.9	78.0	28.8	30.7	55.8	30.2
<b>Ours</b>							
AgentBench-KG agent w/ our tools	<u>98.5</u>	84.4	85.9	97.6	91.7	91.6	68.8
AutoGPT w/ our tools	<b>99.0</b>	<u>97.0</u>	<b>98.5</b>	<u>99.0</u>	<u>99.0</u>	<u>98.5</u>	75.1
AutoGen w/ our tools	<b>99.0</b>	<b>99.0</b>	<u>98.0</u>	<u>99.0</u>	<b>100.0</b>	<b>99.0</b>	<u>77.1</u>
<b>BKGAgent (ours)</b>	89.8	94.1	95.1	<b>100.0</b>	95.1	94.8	<b>78.0</b>

AutoGen (Wu et al., 2023a) and AutoGPT<sup>2</sup>, along with their three improved versions, as well as our BKGAgent, for comparison, as shown in Table 8<sup>3</sup>. Both the final results and process are considered for a more robust evaluation. Since the ground truth is either “support” or “refute”, we use Exact Match (EM) as the metric for the final result. For the process, we employ Qwen2-72B<sup>4</sup> to score based on five criteria: (1) **Understanding**: whether the agent clearly understood the task and the purpose of the given tool. (2) **Reasoning**: whether the agent arrived at the final answer through sufficient evidence and reasoning, rather than simply providing random answers or guessing. (3) **Efficiency**: whether the agent efficiently solved the problem without unnecessary discussion on unrelated topics. (4) **KG Process**: whether the agent queried the knowledge graph during the task. (5) **Information Retrieval**: whether the agent retrieved information from external knowledge sources in some way during the check. To align the judgments made by LLMs closely with those of humans, we collect 10 agent histories along with human-annotated scores (on a 5-point scale) and prompt the LLM to produce scores that closely resemble human ratings. We take EM as the main metric, while process scores serve as supplementary metrics.

**Agent Comparison.** Table 8 shows that BKGAgent outperforms the other agents. KG-Agent achieved an accuracy of only 56.6%, roughly equivalent to random guessing. This aligns with our expectations, as while it can accurately query information from the knowledge graph, it lacks access to reliable external knowledge sources for verification, leading to hallucinations in the large model’s guesses. Notably, the final accuracy of Vanilla AutoGen and AutoGPT is quite low, at just over 30%. This underscores the importance of integrating general capabilities with specialized tools to enhance agent performance. Their performance suffers because they are general frameworks that rely on some general capabilities like programming and web searches, which are not robust enough, often resulting in execution failures due to poor code quality. Consequently, they cannot provide answers within the limited interaction turns. Thus, we improved KG-Agent, AutoGen, and AutoGPT by equipping them with tools including KG querying and RAG. We also designed prompts to teach them how to utilize these tools. As a result, KG-Agent w/ our tools, AutoGen w/ our tools, and AutoGPT w/ our tools demonstrate significant improvements, highlighting that the integration of general capabilities with specialized tools enhances the robustness of agent performance.

**Case Study of BKGAgent.** While the behavior of the assistant agents in BKGAgent can be modified by the leader’s instruction, the leader itself lacks action-related feedback from others, meaning that a bad decision made by the leader may lead to a catastrophe. We found four common error cases induced by the leader, as shown in Figure 6. Among these cases, the leader either fails to give effective instructions to team members, becomes trapped in repeated self-talks, or attempts to perform the tasks that are meant for the assistants. Thus, we can conclude that the performance of the leader agent significantly impacts team behavior. Furthermore, compared to the other agent frameworks in Table 8, BKGAgent employs the maximum number of three agents. However, in our additional case studies (see Appendix D.3), we found that an increase in the number of agents also leads to higher communication costs and a decrease in fault tolerance. They frequently experience task failures due to being lost in conversation.

<sup>2</sup><https://news.agpt.co/>

<sup>3</sup>Other agent frameworks that can perform information retrieval but cannot query knowledge graphs are unable to execute the task from the outset, resulting in scores of zero; therefore, we omit their evaluation.

<sup>4</sup><https://huggingface.co/Qwen/Qwen2-72B-Instruct>



486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

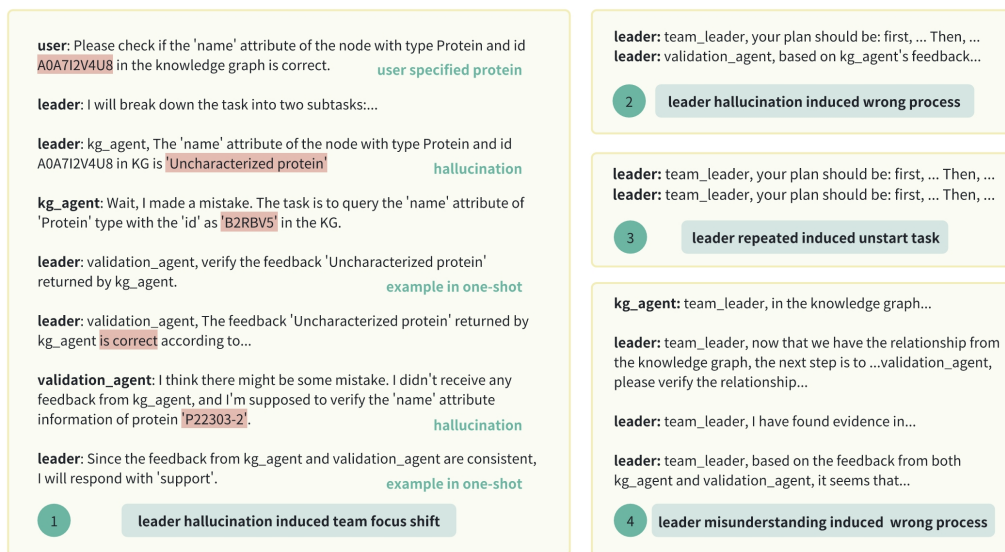


Figure 6: **Error analysis.** Here, we show a failure case due to a leader’s various mistakes: the hallucination of the leader misleading the later task or using the wrong process, the leader producing unnecessary repeated tasks and misunderstanding leads to the wrong process.

### Impact of Agent Number on System Performance: A

**Further Analysis.** By equipping agents with identical capabilities (e.g., KG querying, database querying, and RAG of literature), we compare the performance of systems with 1, 2, and 3 agents (Cf. Figure 5). **Our BKGAgent achieves a 20% higher recall rate of errors (i.e., the ability to correctly identify errors in the KG) than the second-best system, demonstrating its strong performance. However, we also note that the advantage of using 3 agents in terms of EM is not significant.** While adding agents slightly improves performance through collaboration, it also increases communication costs and complexities, leading to diminishing returns. In contrast, the application of effective algorithms, such as ReACT (Yao et al., 2022), can yield more significant improvements, as evidenced by the comparison between AgentBench-KG Agent w/ our tools and AutoGPT w/ our tools.

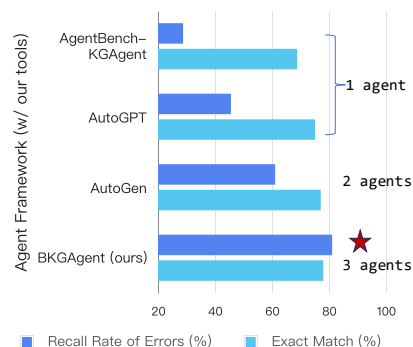


Figure 5: Comparison of the performance of agent frameworks with 1/2/3 agents.

## 5 CONCLUSION

We present BioKGBench, an interactive benchmark that encompasses the KGCheck task with two atomic capabilities for assessment: KGQA and SCV. KGCheck offers agents a valuable scenario for detecting knowledge hallucination within large-scale data, akin to the experience of researchers making discoveries amidst voluminous literature in the real world. We conduct evaluations of these two atomic capabilities across 13 LLMs and select GPT-4o, to construct BKGAgent—a multi-agent system serving as the baseline. Comparisons with existing general and biomedical agents revealed their poor performance due to the absence of certain process capabilities, thereby demonstrating the challenging nature of our benchmark. We expect BioKGBench to serve as a valuable endeavor towards paving the path for biomedical agents to become AI scientists.

**Limitations and Future Work.** In KGCheck, we guide agents to identify knowledge-based errors within the KG by providing them with specific instructions. This process involves atomic-level inspections from single nodes to triples, which agents could potentially implement autonomously. Future work will explore how agents can autonomously conduct real-time error detection in large datasets by leveraging logic rules and prior knowledge.

## REFERENCES

- 540  
541 Hello GPT-4o. <https://openai.com/index/hello-gpt-4o/>.
- 542  
543 Uniprot: the universal protein knowledgebase in 2023. *Nucleic acids research*, 51(D1):D523–D531,  
544 2023.
- 545  
546 Bilal Abu-Salih. Domain-specific knowledge graphs: A survey. *Journal of Network and Computer  
547 Applications*, 185:103076, 2021.
- 548  
549 AI@Meta. Llama 3 model card. 2024. URL [https://github.com/meta-llama/llama3/blob/  
main/MODEL\\_CARD.md](https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md).
- 550  
551 Sakher Khalil Alqaaidi and Krzysztof Kochut. Knowledge graph completion using structural and  
552 textual embeddings. *arXiv preprint arXiv:2404.16206*, 2024.
- 553  
554 Sören Auer, Christian Bizer, Georgi Kobilarov, Jens Lehmann, Richard Cyganiak, and Zachary Ives.  
555 Dbpedia: A nucleus for a web of open data. In *international semantic web conference*, pp. 722–735.  
Springer, 2007.
- 556  
557 Jinheon Baek, Sujay Kumar Jauhar, Silviu Cucerzan, and Sung Ju Hwang. Researchagent: Iterative  
558 research idea generation over scientific literature with large language models. *arXiv preprint  
arXiv:2404.07738*, 2024.
- 559  
560 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge,  
561 Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu,  
562 Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan,  
563 Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin  
564 Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng  
565 Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou,  
566 Xiaohuan Zhou, and Tianhang Zhu. Qwen technical report. *arXiv preprint arXiv:2309.16609*,  
2023.
- 567  
568 DeepSeek-AI Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng,  
569 Honghui Ding, Kai Dong, Qiushi Du, Zhe Fu, Huazuo Gao, Kaige Gao, Wenjun Gao, Ruiqi  
570 Ge, Kang Guan, Daya Guo, Jianzhong Guo, Guangbo Hao, Zhewen Hao, Ying He, Wen-Hui  
571 Hu, Panpan Huang, Erhang Li, Guowei Li, Jiashi Li, Yao Li, Y. K. Li, Wenfeng Liang, Fangyun  
572 Lin, A. X. Liu, Bo Liu, Wen Liu, Xiaodong Liu, Xin Liu, Yiyuan Liu, Haoyu Lu, Shanghao Lu,  
573 Fuli Luo, Shirong Ma, Xiaotao Nie, Tian Pei, Yishi Piao, Junjie Qiu, Hui Qu, Tongzheng Ren,  
574 Zehui Ren, Chong Ruan, Zhangli Sha, Zhihong Shao, Jun-Mei Song, Xuecheng Su, Jingxiang  
575 Sun, Yaofeng Sun, Min Tang, Bing-Li Wang, Peiyi Wang, Shiyu Wang, Yaohui Wang, Yongji  
576 Wang, Tong Wu, Yu Wu, Xin Xie, Zhenda Xie, Ziwei Xie, Yi Xiong, Hanwei Xu, Ronald X Xu,  
577 Yanhong Xu, Dejian Yang, Yu mei You, Shuiping Yu, Xin yuan Yu, Bo Zhang, Haowei Zhang,  
578 Lecong Zhang, Liyue Zhang, Mingchuan Zhang, Minghu Zhang, Wentao Zhang, Yichao Zhang,  
579 Chenggang Zhao, Yao Zhao, Shangyan Zhou, Shunfeng Zhou, Qihao Zhu, and Yuheng Zou.  
580 Deepseek llm: Scaling open-source language models with longtermism. *ArXiv*, abs/2401.02954,  
2024. URL <https://api.semanticscholar.org/CorpusID:266818336>.
- 581  
582 Zhen Bi, Ningyu Zhang, Yida Xue, Yixin Ou, Daxiong Ji, Guozhou Zheng, and Huajun Chen.  
583 OceanGPT: A large language model for ocean science tasks. *arXiv preprint arXiv:2310.02031*,  
2023.
- 584  
585 Daniil A. Boiko, Robert MacKnight, and Gabe Gomes. Emergent autonomous scientific research  
586 capabilities of large language models. *arXiv preprint arXiv:2304.05332*, 2023a.
- 587  
588 Daniil A Boiko, Robert MacKnight, Ben Kline, and Gabe Gomes. Autonomous chemical research  
589 with large language models. *Nature*, 624(7992):570–578, 2023b.
- 590  
591 Kurt Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. Freebase: a collabora-  
592 tively created graph database for structuring human knowledge. In *Proceedings of the 2008 ACM  
SIGMOD international conference on Management of data*, pp. 1247–1250, 2008.
- 593  
Elliot Bolton, David Hall, Michihiro Yasunaga, Tony Lee, Chris Manning, and Percy Liang.  
Biomedlm. <https://crfm.stanford.edu/2022/12/15/biomedlm.html>, 2022.

- 594 Andres M Bran, Sam Cox, Oliver Schilter, Carlo Baldassari, Andrew D White, and Philippe  
595 Schwaller. Chemcrow: Augmenting large-language models with chemistry tools. *arXiv preprint*  
596 *arXiv:2304.05376*, 2023.
- 597
- 598 Hengxing Cai, Xiaochen Cai, Junhan Chang, Sihang Li, Lin Yao, Changxin Wang, Zhifeng Gao,  
599 Yongge Li, Mujie Lin, Shuwen Yang, et al. Sciassess: Benchmarking llm proficiency in scientific  
600 literature analysis. *arXiv preprint arXiv:2403.01976*, 2024.
- 601
- 602 Asli Celikyilmaz, Elizabeth Clark, and Jianfeng Gao. Evaluation of text generation: A survey. *arXiv*  
603 *preprint arXiv:2006.14799*, 2020.
- 604
- 605 Payal Chandak, Kexin Huang, and Marinka Zitnik. Building a knowledge graph to enable precision  
606 medicine. *Scientific Data*, 10(1):67, 2023.
- 607
- 608 Srinivas Niranj Chandrasekaran, Beth A Cimini, Amy Goodale, Lisa Miller, Maria Kost-Alimova,  
609 Nasim Jamali, John G Doench, Briana Fritchman, Adam Skepner, Michelle Melanson, et al. Three  
610 million images and morphological profiles of cells treated with matched chemical and genetic  
611 perturbations. *Nature Methods*, pp. 1–8, 2024.
- 612
- 613 Antje Chang, Ida Schomburg, Sandra Placzek, Lisa Jeske, Marcus Ulbrich, Mei Xiao, Christoph W  
614 Sensen, and Dietmar Schomburg. Brenda in 2015: exciting developments in its 25th year of  
615 existence. *Nucleic acids research*, 43(D1):D439–D446, 2015.
- 616
- 617 Harrison Chase. Langgraph. <https://github.com/langchain-ai/langgraph>, October 2023.  
618 Computer Software.
- 619
- 620 Shengyuan Chen, Yunfeng Cai, Huang Fang, Xiao Huang, and Mingming Sun. Differentiable  
621 neuro-symbolic reasoning on large-scale knowledge graphs. In *Neural Information Processing*  
622 *Systems*, 2023a. URL <https://api.semanticscholar.org/CorpusID:268042075>.
- 623
- 624 Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chen Qian, Chi-Min Chan, Yujia  
625 Qin, Yaxi Lu, Ruobing Xie, et al. Agentverse: Facilitating multi-agent collaboration and exploring  
626 emergent behaviors in agents. *arXiv preprint arXiv:2308.10848*, 2023b.
- 627
- 628 Gene Ontology Consortium. Expansion of the gene ontology knowledgebase and resources. *Nucleic*  
629 *acids research*, 45(D1):D331–D338, 2017.
- 630
- 631 Jean-Michel Dalle and Paul A David. Simcode: agent-based simulation modelling of open-source  
632 software development. *EconWPA, Industrial Organization*, 2004.
- 633
- 634 Noemi Del Toro, Anjali Shrivastava, Eliot Ragueneau, Birgit Meldal, Colin Combe, Elisabet Barrera,  
635 Livia Peretto, Karyn How, Prashansa Ratan, Gautam Shirodkar, et al. The intact database: efficient  
636 access to fine-grained molecular interaction data. *Nucleic acids research*, 50(D1):D648–D653,  
637 2022.
- 638
- 639 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale  
640 hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,  
641 pp. 248–255. Ieee, 2009.
- 642
- 643 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep  
644 bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of*  
645 *the North American Chapter of the Association for Computational Linguistics: Human Language*  
646 *Technologies*, pp. 4171–4186, 2019.
- 647
- 648 Alexis Drogoul and Jacques Ferber. Multi-agent simulation as a tool for studying emergent processes  
649 in societies. In *Simulating societies*, pp. 127–142. Routledge, 2018.
- 650
- 651 Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. GLM:  
652 general language model pretraining with autoregressive blank infilling. pp. 320–335, 2022.
- 653
- 654 Carl N Edwards, Aakanksha Naik, Tushar Khot, Martin D Burke, Heng Ji, and Tom Hope. Synergpt:  
655 In-context learning for personalized drug synergy prediction and drug design. *bioRxiv*, 2023. doi:  
656 10.1101/2023.07.06.547759.

- 648 Wafaa S El-Kassas, Cherif R Salama, Ahmed A Rafea, and Hoda K Mohamed. Automatic text  
649 summarization: A comprehensive survey. *Expert systems with applications*, 165:113679, 2021.  
650
- 651 Antonio Fabregat, Steven Jupe, Lisa Matthews, Konstantinos Sidiropoulos, Marc Gillespie, Phani  
652 Garapati, Robin Haw, Bijay Jassal, Florian Korninger, Bruce May, et al. The reactome pathway  
653 knowledgebase. *Nucleic acids research*, 46(D1):D649–D655, 2018.
- 654 Mahak Gambhir and Vishal Gupta. Recent automatic text summarization techniques: a survey.  
655 *Artificial Intelligence Review*, 47(1):1–66, 2017.  
656
- 657 Shanghua Gao, Ada Fang, Yepeng Huang, Valentina Giunchiglia, Ayush Noori, Jonathan Richard  
658 Schwarz, Yasha Ektefaie, Jovana Kondic, and Marinka Zitnik. Empowering biomedical discovery  
659 with ai agents. *arXiv preprint arXiv:2404.02831*, 2024.
- 660 Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach,  
661 Hal Daumé Iii, and Kate Crawford. Datasheets for datasets. *Communications of the ACM*, 64(12):  
662 86–92, 2021.  
663
- 664 Francesca Grisoni. Chemical language models for de novo drug design: Challenges and opportunities.  
665 *Current Opinion in Structural Biology*, 79:102527, 2023.
- 666 Yu Gu, Sue Kase, Michelle Vanni, Brian Sadler, Percy Liang, Xifeng Yan, and Yu Su. Beyond iid:  
667 three levels of generalization for question answering on knowledge bases. In *Proceedings of the*  
668 *Web Conference 2021*, pp. 3477–3488, 2021.  
669
- 670 Taicheng Guo, Kehan Guo, Zhengwen Liang, Zhichun Guo, Nitesh V Chawla, Olaf Wiest, and  
671 Xiangliang Zhang. What indeed can GPT models do in chemistry? a comprehensive benchmark  
672 on eight tasks. *arXiv preprint arXiv:2305.18365*, 2023.
- 673 Michael Günther, Jackmin Ong, Isabelle Mohr, Alaeddine Abdessalem, Tanguy Abel, Moham-  
674 mad Kalim Akram, Susana Guzman, Georgios Mastrapas, Saba Sturua, Bo Wang, Maximilian  
675 Werk, Nan Wang, and Han Xiao. Jina embeddings 2: 8192-token general-purpose text embeddings  
676 for long documents, 2023.
- 677 Kai He, Rui Mao, Qika Lin, Yucheng Ruan, Xiang Lan, Mengling Feng, and Erik Cambria. A survey  
678 of large language models for healthcare: from data, technology, and applications to accountability  
679 and ethics. *ArXiv*, abs/2310.05694, 2023. URL [https://api.semanticscholar.org/CorpusID:  
680 263829396](https://api.semanticscholar.org/CorpusID:263829396).  
681
- 682 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and  
683 Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint*  
684 *arXiv:2009.03300*, 2020.
- 685 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,  
686 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv*  
687 *preprint arXiv:2103.03874*, 2021.  
688
- 689 Aidan Hogan, Eva Blomqvist, Michael Cochez, Claudia d’Amato, Gerard De Melo, Claudio Gutierrez,  
690 Sabrina Kirrane, José Emilio Labra Gayo, Roberto Navigli, Sebastian Neumaier, et al. Knowledge  
691 graphs. *ACM Computing Surveys (Csur)*, 54(4):1–37, 2021.
- 692 Timothy Jewison, Yilu Su, Fatemeh Miri Disfany, Yongjie Liang, Craig Knox, Adam Maciejewski,  
693 Jenna Poelzer, Jessica Huynh, You Zhou, David Arndt, et al. Smpdb 2.0: big improvements to the  
694 small molecule pathway database. *Nucleic acids research*, 42(D1):D478–D484, 2014.
- 695 Ziwei Ji, Tiezheng Yu, Yan Xu, Nayeon Lee, Etsuko Ishii, and Pascale Fung. Towards mitigating  
696 llm hallucination via self reflection. In *Findings of the Association for Computational Linguistics:*  
697 *EMNLP 2023*, pp. 1827–1843, 2023.  
698
- 699 Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris  
700 Bamford, Devendra Singh Chaplot, Diego de Las Casas, Emma Bou Hanna, Florian Bressand,  
701 Gianna Lengyel, Guillaume Bour, Guillaume Lample, L’elio Renard Lavaud, Lucile Saulnier,  
Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak,

- 702 Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El  
703 Sayed. Mixtral of experts. *ArXiv*, abs/2401.04088, 2024a. URL <https://api.semanticscholar.org/CorpusID:266844877>.  
704
- 705  
706 Jinhao Jiang, Kun Zhou, Zican Dong, Keming Ye, Wayne Xin Zhao, and Ji-Rong Wen. Structgpt:  
707 A general framework for large language model to reason over structured data. *arXiv preprint*  
708 *arXiv:2305.09645*, 2023.
- 709 Jinhao Jiang, Kun Zhou, Wayne Xin Zhao, Yang Song, Chen Zhu, Hengshu Zhu, and Ji-Rong Wen.  
710 Kg-agent: An efficient autonomous agent framework for complex reasoning over knowledge graph.  
711 *arXiv preprint arXiv:2402.11163*, 2024b.
- 712  
713 Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik  
714 Narasimhan. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint*  
715 *arXiv:2310.06770*, 2023.
- 716 Ming Jin, Yifan Zhang, Wei Chen, Kexin Zhang, Yuxuan Liang, Bin Yang, Jindong Wang, Shirui  
717 Pan, and Qingsong Wen. Position paper: What can large language models tell us about time series  
718 analysis. In *International Conference on Machine Learning (ICML 2024)*, 2024a.
- 719 Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William W Cohen, and Xinghua Lu. Pubmedqa: A  
720 dataset for biomedical research question answering. *arXiv preprint arXiv:1909.06146*, 2019.  
721
- 722 Qiao Jin, Yifan Yang, Qingyu Chen, and Zhiyong Lu. Genegpt: Teaching large language models to  
723 use ncbi web apis. 2023.
- 724 Qiao Jin, Yifan Yang, Qingyu Chen, and Zhiyong Lu. Genegpt: Augmenting large language models  
725 with domain tools for improved access to biomedical information. *Bioinformatics*, 40(2):btac075,  
726 2024b.
- 727  
728 Mayank Kejriwal. *Domain-specific knowledge graph construction*. Springer, 2019.
- 729 Yihuai Lan, Zhiqiang Hu, Lei Wang, Yang Wang, Deheng Ye, Peilin Zhao, Ee-Peng Lim, Hui Xiong,  
730 and Hao Wang. Llm-based agent society investigation: Collaboration and confrontation in avalon  
731 gameplay. *arXiv preprint arXiv:2310.14985*, 2023.
- 732  
733 Raymond ST Lee and James NK Liu. ijade web-miner: an intelligent agent framework for internet  
734 shopping. *IEEE Transactions on Knowledge and Data Engineering*, 16(4):461–473, 2004.
- 735 Tianhao Li, Sandesh Shetty, Advait Kamath, Ajay Jaiswal, Xianqian Jiang, Ying Ding, and Yejin  
736 Kim. Cancergpt: Few-shot drug pair synergy prediction using large pre-trained language models.  
737 *arXiv preprint arXiv:2304.10946*, 2023.
- 738  
739 Yutong Li, Lu Chen, Aiwei Liu, Kai Yu, and Lijie Wen. Chatcrite: Llm agent with human workflow  
740 guidance for comparative literature summary. *arXiv preprint arXiv:2403.02574*, 2024.
- 741 Youwei Liang, Ruiyi Zhang, Li Zhang, and Pengtao Xie. Drugchat: Towards enabling chatgpt-like  
742 capabilities on drug molecule graphs. 2023.
- 743  
744 Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig.  
745 Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language  
746 processing. *ACM Computing Surveys*, 55(9):1–35, 2023a.
- 747 Shengchao Liu, Jiong Xiao Wang, Yijin Yang, Chengpeng Wang, Ling Liu, Hongyu Guo, and Chaowei  
748 Xiao. Chatgpt-powered conversational drug editing using retrieval and domain feedback. *arXiv*  
749 *preprint arXiv:2305.18090*, 2023b.
- 750 Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding,  
751 Kaiwen Men, Kejuan Yang, et al. Agentbench: Evaluating llms as agents. *arXiv preprint*  
752 *arXiv:2308.03688*, 2023c.
- 753  
754 Prisca Lo Surdo, Marta Iannuccelli, Silvia Contino, Luisa Castagnoli, Luana Licata, Gianni Cesareni,  
755 and Livia Perfetto. Signor 3.0, the signaling network open resource 3.0: 2022 update. *Nucleic*  
*acids research*, 51(D1):D631–D637, 2023.



- 756 Gerhard Mayer, Luisa Montecchi-Palazzi, David Ovelheiro, Andrew R Jones, Pierre-Alain Binz,  
757 Eric W Deutsch, Matthew Chambers, Marius Kallhardt, Fredrik Levander, James Shofstahl, et al.  
758 The hupo proteomics standards initiative-mass spectrometry controlled vocabulary. *Database*,  
759 2013:bat009, 2013.
- 760  
761 Lars-Peter Meyer, Claus Stadler, Johannes Frey, Norman Radtke, Kurt Junghanns, Roy Meissner,  
762 Gordian Dziwis, Kirill Bulert, and Michael Martin. Llm-assisted knowledge graph engineering:  
763 Experiments with chatgpt. In *Working conference on Artificial Intelligence Development for a*  
764 *Resilient and Sustainable Tomorrow*, pp. 103–115. Springer Fachmedien Wiesbaden Wiesbaden,  
765 2023.
- 766 OpenAI. Gpt-4 technical report. 2023a. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:257532815)  
767 257532815.
- 768 OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023b.
- 769  
770 Sandra Orchard, Mais Ammari, Bruno Aranda, Lionel Breuza, Leonardo Briganti, Fiona Broackes-  
771 Carter, Nancy H Campbell, Gayatri Chavali, Carol Chen, Noemi Del-Toro, et al. The mintact  
772 project—intact as a common curation platform for 11 molecular interaction databases. *Nucleic*  
773 *acids research*, 42(D1):D358–D363, 2014.
- 774 Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. Medmcqa: A large-scale  
775 multi-subject multi-choice dataset for medical domain question answering. In *Conference on*  
776 *health, inference, and learning*, pp. 248–260. PMLR, 2022.
- 777  
778 Oana Palasca, Alberto Santos, Christian Stolte, Jan Gorodkin, and Lars Juhl Jensen. Tissues 2.0: an  
779 integrative web resource on mammalian tissue expression. *Database*, 2018:bay003, 2018.
- 780 Shirui Pan, Linhao Luo, Yufei Wang, Chen Chen, Jiapu Wang, and Xindong Wu. Unifying large  
781 language models and knowledge graphs: A roadmap. *IEEE Transactions on Knowledge and Data*  
782 *Engineering*, 2024.
- 783  
784 Janet Piñero, Juan Manuel Ramírez-Anguita, Josep Saüch-Pitarch, Francesco Ronzano, Emilio  
785 Centeno, Ferran Sanz, and Laura I Furlong. The disgenet knowledge platform for disease genomics:  
786 2019 update. *Nucleic acids research*, 48(D1):D845–D855, 2020.
- 787 Sune Pletscher-Frankild, Albert Pallejà, Kalliopi Tsafou, Janos X Binder, and Lars Juhl Jensen.  
788 Diseases: Text mining and data integration of disease–gene associations. *Methods*, 74:83–89,  
789 2015.
- 790 Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language  
791 understanding by generative pre-training. 2018.
- 792  
793 Igor Rodchenkov, Ozgun Babur, Augustin Luna, Bulent Arman Aksoy, Jeffrey V Wong, Dylan Fong,  
794 Max Franz, Metin Can Siper, Manfred Cheung, Michael Wrana, et al. Pathway commons 2019  
795 update: integration, analysis and exploration of pathway data. *Nucleic acids research*, 48(D1):  
796 D489–D497, 2020.
- 797 Yusuf Roohani, Jian Vora, Qian Huang, Zachary Steinhart, Alexander Marson, Percy Liang, and Jure  
798 Leskovec. Biodiscoveryagent: An ai agent for designing genetic perturbation experiments. *arXiv*  
799 *preprint arXiv:2405.17631*, 2024.
- 800  
801 Alberto Santos, Ana R Colaço, Annelaura B Nielsen, Lili Niu, Maximilian Strauss, Philipp E Geyer,  
802 Fabian Coscia, Nicolai J Wewer Albrechtsen, Filip Mundt, Lars Juhl Jensen, et al. A knowledge  
803 graph to interpret clinical proteomics data. *Nature biotechnology*, 40(5):692–702, 2022.
- 804 Samuel Schmidgall, Rojin Ziaei, Carl Harris, Eduardo Reis, Jeffrey Jopling, and Michael Moor.  
805 Agentclinic: a multimodal agent benchmark to evaluate ai in simulated clinical environments.  
806 *arXiv preprint arXiv:2405.07960*, 2024.
- 807  
808 Lynn M Schriml, Elvira Mitra, James Munro, Becky Tauber, Mike Schor, Lance Nickle, Victor  
809 Felix, Linda Jeng, Cynthia Bearer, Richard Lichenstein, et al. Human disease ontology 2018 update:  
classification, content and workflow expansion. *Nucleic acids research*, 47(D1):D955–D962, 2019.

- 810 Ruth L Seal, Bryony Braschi, Kristian Gray, Tamsin EM Jones, Susan Tweedie, Liora Haim-  
811 Vilmovsky, and Elspeth A Bruford. Genenames. org: the hgnc resources in 2023. *Nucleic*  
812 *Acids Research*, 51(D1):D1003–D1009, 2023.
- 813 Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. Hugginggpt:  
814 Solving ai tasks with chatgpt and its friends in hugging face, 2023.
- 815 K. Singhal, Shekoofeh Azizi, Tao Tu, Said Mahdavi, Jason Wei, Hyung Won Chung, Nathan  
816 Scales, Ajay Kumar Tanwani, Heather J. Cole-Lewis, Stephen J. Pfohl, P A Payne, Martin G.  
817 Seneviratne, Paul Gamble, Chris Kelly, Nathaneal Scharli, Aakanksha Chowdhery, P. A. Mansfield,  
818 Blaise Agüera y Arcas, Dale R. Webster, Greg S. Corrado, Yossi Matias, Katherine Hui-Ling  
819 Chou, Juraj Gottweis, Nenad Tomasev, Yun Liu, Alvin Rajkomar, Joëlle K. Barral, Christopher  
820 Semturs, Alan Karthikesalingam, and Vivek Natarajan. Large language models encode clinical  
821 knowledge. *Nature*, 620:172 – 180, 2022. URL [https://api.semanticscholar.org/CorpusID:  
822 255124952](https://api.semanticscholar.org/CorpusID:255124952).
- 823 K. Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Le Hou, Kevin Clark, Stephen R.  
824 Pfohl, Heather J. Cole-Lewis, Darlene Neal, Mike Schaekermann, Amy Wang, Mohamed Amin,  
825 S. Lachgar, P. A. Mansfield, Sushant Prakash, Bradley Green, Ewa Dominowska, Blaise Agüera  
826 y Arcas, Nenad Tomasev, Yun Liu, Renee C Wong, Christopher Semturs, Seyedeh Sara Mah-  
827 davi, Joëlle K. Barral, Dale R. Webster, Greg S Corrado, Yossi Matias, Shekoofeh Azizi, Alan  
828 Karthikesalingam, and Vivek Natarajan. Towards expert-level medical question answering with  
829 large language models. *arXiv preprint arXiv:2305.09617*, 2023.
- 830 Hannah Snyder. Literature review as a research methodology: An overview and guidelines. *Journal*  
831 *of business research*, 104:333–339, 2019.
- 832 Peiyang Song, Kaiyu Yang, and Anima Anandkumar. Towards large language models as copilots for  
833 theorem proving in Lean. *arXiv preprint arXiv: Arxiv-2404.12534*, 2024.
- 834 Vladimir Sotnikov and Anastasiia Chaikova. Language models for multimessenger astronomy.  
835 *Galaxies*, 11(3), 2023. ISSN 2075-4434. doi: 10.3390/galaxies11030063. URL [https://www.  
836 mdpi.com/2075-4434/11/3/63](https://www.mdpi.com/2075-4434/11/3/63).
- 837 Neil Joseph Steiner. *Autonomous computing systems*. PhD thesis, Virginia Tech, 2008.
- 838 Jiashuo Sun, Chengjin Xu, Lumingyuan Tang, Saizhuo Wang, Chen Lin, Yeyun Gong, Heung-Yeung  
839 Shum, and Jian Guo. Think-on-graph: Deep and responsible reasoning of large language model  
840 with knowledge graph. *arXiv preprint arXiv:2307.07697*, 2023.
- 841 Damian Szklarczyk, Rebecca Kirsch, Mikaela Koutrouli, Katerina Nastou, Farrokh Mehryary, Radja  
842 Hachilif, Annika L Gable, Tao Fang, Nadezhda T Doncheva, Sampo Pyysalo, et al. The string  
843 database in 2023: protein–protein association networks and functional enrichment analyses for any  
844 sequenced genome of interest. *Nucleic acids research*, 51(D1):D638–D646, 2023.
- 845 Alon Talmor and Jonathan Berant. The web as a knowledge-base for answering complex questions.  
846 *arXiv preprint arXiv:1803.06643*, 2018.
- 847 Jiabin Tang, Yuhao Yang, Wei Wei, Lei Shi, Lixin Su, Suqi Cheng, Dawei Yin, and Chao Huang.  
848 Graphgpt: Graph instruction tuning for large language models. *ArXiv*, abs/2310.13023, 2023a.  
849 URL <https://api.semanticscholar.org/CorpusID:264405943>.
- 850 Xiangru Tang, Anni Zou, Zhuosheng Zhang, Yilun Zhao, Xingyao Zhang, Arman Cohan, and Mark  
851 Gerstein. Medagents: Large language models as collaborators for zero-shot medical reasoning.  
852 *arXiv preprint arXiv:2311.10537*, 2023b.
- 853 Xiangru Tang, Anni Zou, Zhuosheng Zhang, Yilun Zhao, Xingyao Zhang, Arman Cohan, and  
854 Mark B. Gerstein. Medagents: Large language models as collaborators for zero-shot medical rea-  
855 soning. *ArXiv*, abs/2311.10537, 2023c. URL [https://api.semanticscholar.org/CorpusID:  
856 265281260](https://api.semanticscholar.org/CorpusID:265281260).
- 857 Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony S. Hartshorn, Elvis  
858 Saravia, Andrew Poulton, Viktor Kerkez, and Robert Stojnic. Galactica: A large language model  
859 for science. *arXiv preprint arXiv:2211.09085*, 2022.
- 860  
861  
862  
863

- 864 XAgent Team. Xagent: An autonomous agent for complex task solving, 2023.  
865
- 866 Julie Thomas, Juliana Utley, Soo-Young Hong, Hunkar Korkmaz, and Gwen Nugent. A review of the  
867 research. *Handbook of Research on STEM Education*, 2020.  
868
- 869 Yu Tian, Xiao Yang, Jingyuan Zhang, Yinpeng Dong, and Hang Su. Evil geniuses: Delving into the  
870 safety of llm-based agents. *arXiv preprint arXiv:2311.11855*, 2023.  
871
- 872 Denny Vrandečić and Markus Krötzsch. Wikidata: a free collaborative knowledgebase. *Communica-*  
873 *tions of the ACM*, 57(10):78–85, 2014.
- 874 David Wadden, Shanchuan Lin, Kyle Lo, Lucy Lu Wang, Madeleine van Zuylen, Arman Cohan, and  
875 Hannaneh Hajishirzi. Fact or fiction: Verifying scientific claims. *arXiv preprint arXiv:2004.14974*,  
876 2020.  
877
- 878 Hanchen Wang, Tianfan Fu, Yuanqi Du, Wenhao Gao, Kexin Huang, Ziming Liu, Payal Chandak,  
879 Shengchao Liu, Peter Van Katwyk, Andreea Deac, Anima Anandkumar, Karianne J. Bergen,  
880 Carla P. Gomes, Shirley Ho, Pushmeet Kohli, Joan Lasenby, Jure Leskovec, Tie-Yan Liu, Arjun K.  
881 Manrai, Debora S. Marks, Bharath Ramsundar, Le Song, Jimeng Sun, Jian Tang, Petar Velickovic,  
882 Max Welling, Linfeng Zhang, Connor W. Coley, Yoshua Bengio, and Marinka Zitnik. Scientific  
883 discovery in the age of artificial intelligence. *Nature*, 620:47–60, 2023a. URL <https://api.semanticscholar.org/CorpusID:260384616>.  
884
- 885 Hanchen Wang, Tianfan Fu, Yuanqi Du, Wenhao Gao, Kexin Huang, Ziming Liu, Payal Chandak,  
886 Shengchao Liu, Peter Van Katwyk, Andreea Deac, et al. Scientific discovery in the age of artificial  
887 intelligence. *Nature*, 620(7972):47–60, 2023b.  
888
- 889 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny  
890 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in*  
891 *neural information processing systems*, 35:24824–24837, 2022.
- 892 Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li,  
893 Li Jiang, Xiaoyun Zhang, and Chi Wang. Autogen: Enabling next-gen llm applications via  
894 multi-agent conversation framework. *arXiv preprint arXiv:2308.08155*, 2023a.  
895
- 896 Tianyu Wu, Shizhu He, Jingping Liu, Siqi Sun, Kang Liu, Qing-Long Han, and Yang Tang. A brief  
897 overview of chatgpt: The history, status quo and potential future development. *IEEE/CAA Journal*  
898 *of Automatica Sinica*, 10(5):1122–1136, 2023b.  
899
- 900 Yihang Xiao, Jinyi Liu, Yan Zheng, Xiaohan Xie, Jianye Hao, Mingzhi Li, Ruitao Wang, Fei Ni,  
901 Yuxiao Li, Jintian Luo, et al. Cellagent: An llm-driven multi-agent framework for automated  
902 single-cell data analysis. *bioRxiv*, pp. 2024–05, 2024.
- 903 Tianbao Xie, Fan Zhou, Zhoujun Cheng, Peng Shi, Luoxuan Weng, Yitao Liu, Toh Jing Hua, Junning  
904 Zhao, Qian Liu, Che Liu, et al. Openagents: An open platform for language agents in the wild.  
905 *arXiv preprint arXiv:2310.10634*, 2023.  
906
- 907 Ai Ming Yang, Bin Xiao, Bingning Wang, Borong Zhang, Ce Bian, Chao Yin, Chenxu Lv, Da Pan,  
908 Dian Wang, Dong Yan, Fan Yang, Fei Deng, Feng Wang, Feng Liu, Guangwei Ai, Guosheng Dong,  
909 Hai Zhao, Hang Xu, Hao-Lun Sun, Hongda Zhang, Hui Liu, Jiaming Ji, Jian Xie, Juntao Dai,  
910 Kuncheng Fang, Lei Su, Liang Song, Lifeng Liu, Liyun Ru, Luyao Ma, Mang Wang, Mickel Liu,  
911 MingAn Lin, Nuolan Nie, Pei Guo, Ruiyang Sun, Zhang Tao, Tianpeng Li, Tianyu Li, Wei Cheng,  
912 Weipeng Chen, Xiangrong Zeng, Xiaochuan Wang, Xiaoxi Chen, Xin Men, Xin Yu, Xuehai Pan,  
913 Yan-Bin Shen, Yiding Wang, Yiyu Li, Youxin Jiang, Yuchen Gao, Yupeng Zhang, Zenan Zhou,  
914 and Zhiying Wu. Baichuan 2: Open large-scale language models. *ArXiv*, abs/2309.10305, 2023a.  
915 URL <https://api.semanticscholar.org/CorpusID:261951743>.
- 916 Kaiyu Yang, Aidan Swope, Alex Gu, Rahul Chalamala, Peiyang Song, Shixing Yu, Saad Godil, Ryan  
917 Prenger, and Anima Anandkumar. LeanDojo: Theorem proving with retrieval-augmented language  
models. In *Neural Information Processing Systems (NeurIPS)*, 2023b.

- 918 Lin F. Yang, Hongyang Chen, Zhao Li, Xiao Ding, and Xindong Wu. Give us the facts:  
919 Enhancing large language models with knowledge graphs for fact-aware language model-  
920 ing. *IEEE Transactions on Knowledge and Data Engineering*, 36:3091–3110, 2023c. URL  
921 <https://api.semanticscholar.org/CorpusID:259203671>.
- 922 Jia-Yu Yao, Kun-Peng Ning, Zhen-Hui Liu, Mu-Nan Ning, and Li Yuan. Llm lies: Hallucinations are  
923 not bugs, but features as adversarial examples. *arXiv preprint arXiv:2310.01469*, 2023.
- 924 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao.  
925 React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*,  
926 2022.
- 927 Wen-tau Yih, Matthew Richardson, Christopher Meek, Ming-Wei Chang, and Jina Suh. The value  
928 of semantic parse labeling for knowledge base question answering. In *Proceedings of the 54th*  
929 *Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pp.  
930 201–206, 2016.
- 931 Yoheinakajima. Yoheinakajima/babyagi. URL <https://github.com/yoheinakajima/babyagi>.
- 932 Wenhao Yu, Chenguang Zhu, Zaitang Li, Zhiting Hu, Qingyun Wang, Heng Ji, and Meng Jiang. A  
933 survey of knowledge-enhanced text generation. *ACM Computing Surveys*, 54(11s):1–38, 2022.
- 934 Cyril Zakka, Rohan Shad, Akash Chaurasia, Alex R Dalal, Jennifer L Kim, Michael Moor, Robyn  
935 Fong, Curran Phillips, Kevin Alexander, Euan Ashley, et al. Almanac—retrieval-augmented  
936 language models for clinical medicine. *NEJM AI*, 1(2):A10a2300068, 2024.
- 937 Lei Zhang, Yuge Zhang, Kan Ren, Dongsheng Li, and Yuqing Yang. Mlcopilot: Unleashing the power  
938 of large language models in solving machine learning tasks. *arXiv preprint arXiv:2304.14979*,  
939 2023.
- 940 Xiao Zhang, Ji Wu, Zhiyang He, Xien Liu, and Ying Su. Medical exam question answering with  
941 large-scale reading comprehension. In *Proceedings of the AAAI conference on artificial intelligence*,  
942 volume 32, 2018.
- 943 Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min,  
944 Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Z. Chen, Jinhao  
945 Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jianyun Nie, and Ji rong  
946 Wen. A survey of large language models. *ArXiv*, abs/2303.18223, 2023. URL <https://api.semanticscholar.org/CorpusID:257900969>.
- 947 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,  
948 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and  
949 chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623, 2023.
- 950 Shuangjia Zheng, Jiahua Rao, Ying Song, Jixian Zhang, Xianglu Xiao, Evandro Fei Fang, Yuedong  
951 Yang, and Zhangming Niu. Pharmkg: a dedicated knowledge graph benchmark for biomedical data  
952 mining. *Briefings in bioinformatics*, 22(4):bbaa344, 2021.
- 953 Li Zhenzhu, Zhang Jingfeng, Zhou Wei, Zheng Jianjun, and Xia Yinshui. Gpt-agents based on  
954 medical guidelines can improve the responsiveness and explainability of outcomes for traumatic  
955 brain injury rehabilitation. *Scientific Reports*, 14(1):7626, 2024.
- 956 Ruiqi Zhong, Peter Zhang, Steve Li, Jinwoo Ahn, Dan Klein, and Jacob Steinhardt. Goal driven  
957 discovery of distributional differences via language descriptions. *arXiv preprint arXiv:2302.14233*,  
958 2023.
- 959 Hongjian Zhou, Boyang Gu, Xinyu Zou, Yiru Li, Sam S. Chen, Peilin Zhou, Junling Liu, Yining  
960 Hua, Chengfeng Mao, Xian Wu, Zheng Li, and Fenglin Liu. A survey of large language models  
961 in medicine: Progress, application, and challenge. *ArXiv*, abs/2311.05112, 2023. URL <https://api.semanticscholar.org/CorpusID:265067484>.
- 962 Banghua Zhu, Evan Frick, Tianhao Wu, Hanlin Zhu, and Jiantao Jiao. Starling-7b: Improving llm  
963 helpfulness & harmlessness with rlaif, November 2023.

## A DATASHEET

Here, we provide a detailed description of our benchmark dataset, following the guidelines of the “Datasheet for Datasets” (Geburu et al., 2021).

### A.1 MOTIVATION

Our benchmark dataset was created to address the lack of benchmarks for evaluating biomedical agents from the perspective of an “AI scientists”. In (Gao et al., 2024), it is stated that “AI scientists” can be realized as AI agents supported by humans, LLMs, ML models, and other tools like experimental platforms that cooperate to solve complex tasks. However, the current evaluation methods for biomedical agents remain unexplored, limited to simple question-answering tasks, which not only fail to avoid the hallucination problem inherent in solely relying on LLMs but also do not assess agents’ abilities to utilize external tools and knowledge bases. Our proposed benchmark fills this gap by designing tasks ranging from easy to hard, based on two atomic capabilities: tool-query and memory-RAG. These tasks evaluate the agents’ ability to leverage external support, including external knowledge and tools, when handling large-scale and multi-modal data. Moreover, Our data collected for KGCheck, the most challenging task, provides scenarios for agents to comprehend knowledge from heterogeneous sources and make discoveries.

The conception and construction of this dataset were jointly completed by the biomedical experts and AI researchers listed in the author list.

### A.2 COMPOSITION

We provide the necessary data for constructing the knowledge graph, literature for RAG, as well as the development and test data for KGQA, SCV, and KGCheck, where knowledge graph and literature are external knowledge sources provided for agent.

The knowledge graph is derived from a subset of Clinical Knowledge Graph (CKG) (Santos et al., 2022). We specifically retain twelve key node types to ensure there is no loss of generality: Protein, Biological process (BP), Molecular function (MF), Cellular component (CC), Amino acid sequence, Tissue, Protein structure, Pathway, Modified protein, Modification, Disease, and Gene. The statistics of the triples in our knowledge graph are presented in Table 9. Detailed information stored in our knowledge graph is shown in Table 10.

Table 9: The data statistics of our knowledge graph drawn from CKG.

Head Node	Tail Node	Relation	Number
Protein	Protein_structure	HAS_STRUCTURE	271,512
	Amino_acid_sequence	HAS_SEQUENCE	20,598
	Cellular_component	ASSOCIATED_WITH	3,796,383
	Tissue	ASSOCIATED_WITH	7,117,321
	Disease	ASSOCIATED_WITH	5,882,437
	Molecular_function	ASSOCIATED_WITH	85,013
	Biological_process	ASSOCIATED_WITH	153,219
	Protein	ACTS_ON	985,376
	Pathway	ANNOTATED_IN_PATHWAY	357,739
	Protein	CURATED_INTERACTS_WITH	3,448
	Modified_protein	HAS_MODIFIED_SITE	4,498
Disease	Disease	HAS_PARENT	16,058
Modified_protein	Protein	IS_SUBSTRATE_OF	6,633
	Modification	HAS_MODIFICATION	4,559
Gene	Protein	TRANSLATED_INTO	179,854
Biological_process	Biological_process	HAS_PARENT	49,081
Molecular_function	Molecular_function	HAS_PARENT	13,659

Besides the knowledge graph, literature also serves as an external source of knowledge. We provide a corpus of 5,664 abstracts (under ODC-By 1.0) for SCV and 51 full papers for KGCheck. The 5,664



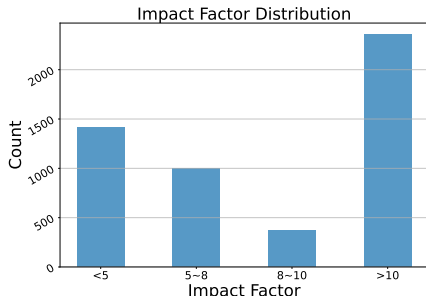
Table 10: Details of the information stored in the nodes of our knowledge graph.

Entity Type	Content	Example
Protein	name, id, accession, synonyms	{ 'name': 'PLEKHG6', 'id': 'Q3KR16', 'accession': 'PKHG6_HUMAN', 'synonyms': ['PKHG6_HUMAN', 'PLEKHG6', '9606.ENSP00000380185', 'ENSG00000008323'], 'taxid': '9606' }
Disease	name, description, id(DOID), type, synonyms	{ 'synonyms': ['sulfamethoxazole allergy', 'SMX allergy', 'SMZ allergy', 'sulphamethoxazole allergy'], 'name': 'sulfamethoxazole allergy', 'description': 'A drug allergy that has_allergic_trigger sulfamethoxazole. [url:https://www.ncbi.nlm.nih.gov/pubmed/7602118]', 'id': 'DOID:0040016', 'type': '-26' }
Protein structure	link, id, source	{ 'link': 'http://www.rcsb.org/structure/6XWD', 'id': '6XWD', 'source': 'Uniprot' }
Amino acid sequence	sequence, header, source, id, size	{ 'sequence': 'LRGAAGRLGGGLLVL', 'size': '15', 'header': 'trIA0A0A0MTA2 A0A0A0MTA2_HUMAN', 'source': 'UniProt', 'id': 'A0A0A0MTA2' }
Cellular component	name, description, id, type, synonyms	{ 'name': 'Golgi membrane', 'description': 'The lipid bilayer surrounding any of the compartments of the Golgi apparatus. [GOC:mah]', 'id': 'GO:0000139', 'type': '-21', 'synonyms': ['Golgi membrane'] }
Molecular function	name, description, id(GO), type, synonyms	{ 'name': 'polymeric immunoglobulin binding', 'description': 'Interacting selectively and non-covalently with a J-chain-containing polymeric immunoglobulin of the IgA or IgM isotypes. [GOC:add, ISBN:0781735149]', 'id': 'GO:0001790', 'type': '-21', 'synonyms': ['polymeric immunoglobulin binding'] }
Biological process	name, description, id(GO), type, synonyms	{ 'synonyms': ['mitochondrion inheritance'], 'name': 'mitochondrion inheritance', 'description': 'The distribution of mitochondria, including the mitochondrial genome, into daughter cells after mitosis or meiosis, mediated by interactions between mitochondria and the cytoskeleton. [GOC:mcc, PMID:10873824, PMID:11389764]', 'id': 'GO:0000001', 'type': '-21' }
Pathway	name, description, linkout, id, source	{ 'name': 'Antigen processing: Ubiquitination & Proteasome degradation', 'description': 'Antigen processing: Ubiquitination & Proteasome degradation', 'linkout': 'https://reactome.org/PathwayBrowser/#/R-HSA-983168', 'id': 'R-HSA-983168', 'source': 'Reactome' }
Tissue	name, description, id, type, synonyms	{ 'name': 'stratum basale', 'description': 'The deepest layer, as of the epidermis or the endometrium. In the epidermis it is a single layer of cells. In the endometrium it provides the regenerative tissue after menstrual loss of the functional layer. [Dorlands_Medical_Dictionary:MerckMedicus]', 'id': 'BTO:0004680', 'type': '-25', 'synonyms': ['stratum basale', 'basal layer'] }
Modified protein	sequence_window, protein, position, source, id, residue	{ 'sequence_window': 'MEPAPARsPRPQQDP', 'protein': 'P29590', 'position': '8', 'source': 'SIGNOR', 'id': 'P29590_S8-p', 'residue': 'S' }
Modification	synonyms, name, description, id, type	{ 'synonyms': ['Unimod', 'Source: "none"'], 'name': 'Unimod', 'description': 'Entry from Unimod. [PubMed:18688235]', 'id': 'MOD:00003', 'type': '-41' }
Gene	taxid, synonyms, name, id, family	{ 'taxid': '9606', 'synonyms': ['54843', 'ENSG00000137501', 'OTTHUMG00000166977', 'uc010rti.4', 'AJ303364'], 'name': 'synaptotagmin like 2', 'id': 'SYTL2', 'family': "'Protein phosphatase 1 regulatory subunits Synaptotagmin like tandem C2 proteins'" }

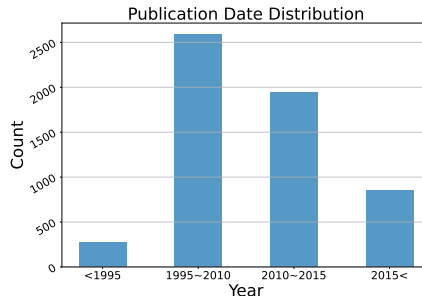
abstracts are sourced from existing datasets SciFact (Wadden et al., 2020) (under CC BY 4.0) and PubMedQA (Jin et al., 2019) (under MIT license), while the 51 full papers, all of which are open access, were selected by experts based on entries in the IntAct (Orchard et al., 2014) database and CKG. Table 11 summarizes the sources of the abstracts, and Figure 7 describes the literature with more details. Table 12 summarizes the sources of the 51 full papers, and Figure 8 provides more details about the literature.

Table 11: The 5,664 papers come from 1,484 journals. Due to space limitations, we only list the names of journals with an IF greater than 70 and use ‘others’ to represent journals with an IF less than 70.

Journal	Count
Nature reviews. Microbiology	3
CA: a cancer journal for clinicians	3
The Lancet. Infectious diseases	3
Nature reviews. Drug discovery	5
Nature reviews. Molecular cell biology	14
Nature reviews. Immunology	21
Lancet (London, England)	46
The New England journal of medicine	46
BMJ (Clinical research ed.)	90
JAMA	113
Nature medicine	138
others	5182
<b>Total</b>	<b>5664</b>



(a) Distribution of literature IF (Impact Factor)



(b) Distribution of literature publication dates

Figure 7: Specific information of the 5,664 abstracts

For evaluation, we carefully collected 698, 1385, and 225 instances for KGQA, SCV, and KGCheck respectively. These datasets are split into development (dev) and test sets at an approximate ratio of 1:10. The dev data is intended for users to debug and fine-tune their evaluation code, while the test data is reserved for the final assessment. Each instance includes both input and output (ground truth answer or label) pairs, with additional information to make the data easier to understand. The dataset for SCV is reconstructed from well-known existing datasets SciFact and PubMedQA, while the rest is self-contained. The dataset represents a carefully selected sample of instances from a larger set, ensuring a comprehensive and representative coverage of the key aspects.

### A.3 COLLECTION PROCESS

Biomedical experts and AI researchers listed in the author list were involved in the data collection process. The collection process for different tasks varies:

**KGQA.** The collection process can be summarized into two steps: manually constructing question templates and automatically generating questions:

Table 12: Journal distribution of the 51 full papers

Journal	Count
Brain research	1
Molecular & cellular proteomics : MCP	1
Aging cell	1
Cell reports	1
PloS one	1
Genes to cells : devoted to molecular & cellular mechanisms	1
EMBO reports	1
IUBMB life	1
The Journal of allergy and clinical immunology	1
The EMBO journal	1
Open biology	1
Nature communications	1
Pigment cell & melanoma research	1
Molecular biology of the cell	1
Mobile DNA	1
Journal of molecular biology	1
Nutrients	1
Biological research for nursing	1
Genes & development	1
Developmental cell	1
Bone	1
Cancers	1
Animals : an open access journal from MDPI	1
Nucleic acids research	2
Molecular cell	2
Scientific reports	3
Proceedings of the National Academy of Sciences of the United States of America	3
Molecular and cellular biology	4
Cell	4
The Journal of biological chemistry	10
Total	51

- **Workflow for the Handcrafted Question Templates:** The process commenced by selecting specific biomedical research fields and identifying relevant entity types and relationships from our knowledge graph. We defined various types of natural language questions, including one-hop questions, multi-hop questions, and conjunction questions (involving multiple entities). For each question, we created corresponding queries in both human-readable and machine-readable formats. These questions and queries, along with their associated metadata, such as question type and query structure, underwent rigorous peer reviews to ensure syntactic and semantic correctness.
- **Workflow for the Auto-generated Questions and Answers:** In the expansion of our benchmark, we initiated the process with the creation of auto-generated question templates. For instance, we used handcrafted question templates like "Which pathway are the proteins <Protein1> and <Protein2> both annotated in?" and then scoured our knowledge graph for data that fit the criteria to formulate both questions and answers, thereby augmenting the size of our dataset. This comprehensive dataset enables the development of assessing the robustness and accuracy of various LLM agents, providing a comprehensive benchmark that contributes to the advancement of the field with extensive biomedical knowledge.

**SCV.** We combine two high-quality biomedical datasets, PubMedQA and SciFact, into a single dataset for SCV. This results in a corpus consisting of abstracts from 5,664 scholarly articles and a dataset containing 1,385 biomedical claims. [To further ensure consistency, we conducted secondary](#)

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

1201

1202

1203

1204

1205

1206

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1220

1221

1222

1223

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

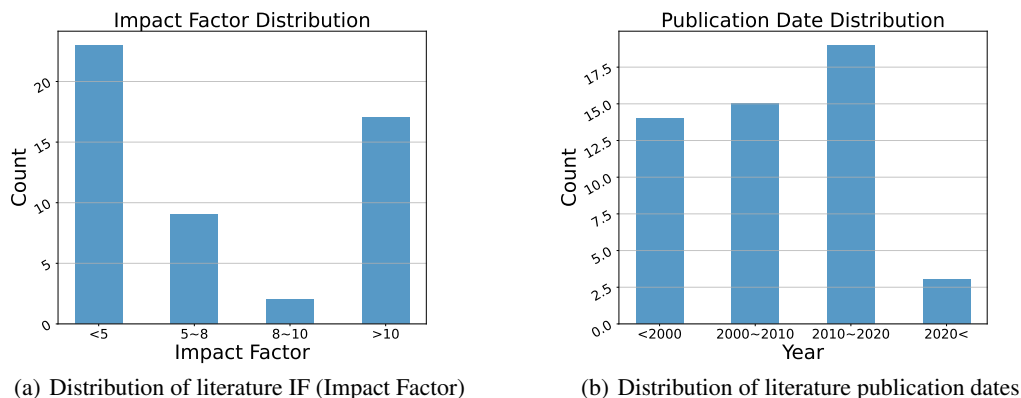


Figure 8: Specific information of the 51 full papers

verification on this expert-annotated dataset using Qwen1.5-72B (Bai et al., 2023), confirming the claims are conflict-free.

**KGCheck.** Considering the characteristics of the knowledge graph, we decompose the approach to checking the knowledge graph into two atomic-level checks: nodes and triples. Further, we subdivide these into whether a node should exist in the knowledge graph, whether the information stored in the node is correct, whether the relationship between two connected nodes in the knowledge graph truly exists, and whether there is a potential relationship between two nodes that are not connected by an edge. To collect this data, we selected well-maintained external knowledge sources such as the UniProt database, the IntAct database, and literature. We cross-verified the information in our knowledge graph with these reliable sources, labeling mutually corroborative data as ‘support’ and data that contradicts the external reliable sources as ‘refute’. Specifically, for the data collection to check nodes, we review some update information from databases, such as entries removed due to errors or entries with updated information. Based on this information, we used Cypher queries to check our knowledge graph and label the data accordingly. For checking triple relationships, we sampled some triples from our knowledge graph where two nodes were either related or unrelated. We then queried the CKG to obtain literature that documents the entities represented by both nodes. We collected the literature annotated in the database and had experts read these documents to label the relationships of the triples in the CKG. As a result, we obtained 225 high-quality annotated data.

#### A.4 USES

The dataset has not been used for any tasks yet. Currently, we have not identified any tasks that are not permitted to use our dataset.

The way we collect question and answer pairs can be referenced to expand more KGQA data, whether on our knowledge graph or new knowledge graphs. Additionally, our approach to collecting data for KGCheck provides insights into identifying errors in these large knowledge graphs, which is very helpful for subsequent error correction and data updates.

## B BREAKDOWN RESULTS

### B.1 KGQA

We conducted a more detailed evaluation of LLMs’ performance on the KGQA task based on the question types: one-hop, multi-hop, and conjunction. The evaluation metrics used were F1 and executability, as shown in Table 13. We find that although API-based commercial LLMs and large-scale open-source models generally perform well on overall metrics, when breaking down the KGQA task by question type, some medium-scale models perform better on certain metrics. For instance, Qwen1.5-14B-Chat exhibits higher executability on more challenging multi-hop and conjunction

Table 13: KGQA Test set (standard) results by question type: one-hop, multi-hop, and conjunction. **Bold/underline** and **red/blue** indicate the best and second in the subgroup and overall.

LLM Type	Models	F1			executability		
		one-hop	multi-hop	conjunction	one-hop	multi-hop	conjunction
API	GPT-4	<b>87.2</b>	<b>73.7</b>	<u>77.4</u>	<b>88.0</b>	<b>90.0</b>	<b>86.9</b>
	GLM-4	76.0	73.0	58.0	82.9	<b>90.0</b>	<u>68.2</u>
OSS (Large)	Qwen1.5-72B-Chat	76.3	<b>73.4</b>	<u>71.4</u>	<b>99.7</b>	94.0	<u>86.9</u>
	Llama-3-70B-Instruct	<b>83.6</b>	<u>72.5</u>	<b>85.1</b>	<u>95.7</u>	<b>98.5</b>	<b>99.1</b>
	DeepSeek-LLM-67B-Chat	<u>80.6</u>	61.8	44.1	88.5	90.5	70.1
OSS (Medium)	Qwen1.5-32B-Chat	<b>67.3</b>	<u>63.2</u>	<u>57.0</u>	<b>87.2</b>	84.5	64.5
	Qwen1.5-14B-Chat	63.7	<b>70.5</b>	<b>65.7</b>	67.5	<b>95.5</b>	<b>87.9</b>
	Baichuan2-13B-Chat	64.9	20.4	9.8	<u>81.8</u>	<u>91.5</u>	<u>66.4</u>
OSS (Small)	Llama-3-8B-Instruct	<b>59.2</b>	<b>66.4</b>	<u>16.5</u>	<b>90.8</b>	<u>66.4</u>	<b>68.2</b>
	Qwen1.5-7B-chat	<u>55.7</u>	<u>32.1</u>	<b>26.4</b>	80.3	<b>79.0</b>	<u>67.3</u>
OSS (MoE)	Mixtral-8x7B-Instruct-v0.1	<b>80.3</b>	<b>68.4</b>	<b>35.9</b>	<u>90.5</u>	<b>91.5</b>	<b>50.5</b>
	Starling-LM-alpha-8x7B-MoE-GPTQ	6.2	25.0	<u>11.7</u>	12.0	<u>57.5</u>	<u>48.6</u>
	Qwen1.5-MoE-A2.7B-Chat	<u>38.2</u>	20.2	9.7	<b>94.4</b>	45.0	40.2

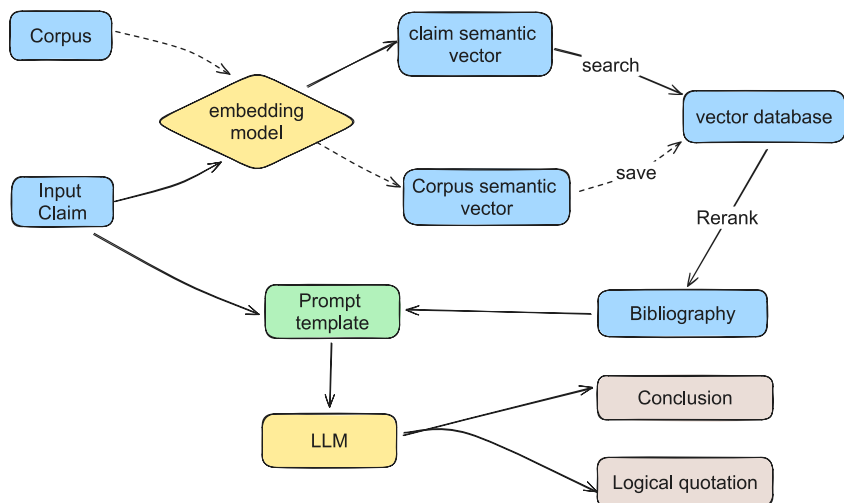


Figure 9: The pipeline of RAG.

types of questions, although its F1 score is not high. In terms of the executability metric, open-source models seem to outperform API-based commercial LLMs. This may be because API-based LLMs are more cautious in determining whether an answer has been obtained, tending to conclude the interaction and return an answer only after confirming its correctness.

## B.2 SCV

As shown in Figure 9, this is the process we followed when performing the SCV task using RAG. In the main text, we observed an interesting phenomenon where expanding the RAG scope improved accuracy. To ensure that this result was not due to the idiosyncratic performance differences of a single model, we conducted the same experiment on another model, as shown in Table 14. It can be observed that the accuracy of both models in the SCV task increased with the expansion of the RAG scope, although the right quotes metric was the lowest across the three settings when performing RAG at the maximum scope. This experimental result further demonstrates that this interesting phenomenon is not due to model-specific characteristics.



Table 14: Supplementary Experiments on the Scope of RAG, where ‘all’ refers to the abstract of 5,664 articles, ‘partial’ denotes the 1,888 abstracts containing ground truth evidence of claims, and ‘match’ corresponds to the abstracts of the ground truth evidence for the claims. **Bold/underline** indicate the best and suboptimal.

Corpus	Qwen1.5-72B-Chat			Llama-3-70B-Instruct		
	accuracy	right quotes	error	accuracy	right quotes	error
all	<b>86.2</b>	82.1	<b>0.1</b>	<b>86.0</b>	86.9	<u>0.2</u>
partial	84.4	88.1	<b>0.1</b>	<u>82.2</u>	<b>90.4</b>	<b>0.1</b>
match	<u>84.3</u>	<b>88.2</b>	<b>0.1</b>	82.1	<u>89.9</u>	<b>0.1</b>

### B.3 KGCHECK

#### B.3.1 SINGLE AGENT

We develop a single agent based on KG-Agent from (Liu et al., 2023c), shown in Figure 10, and evaluate LLMs as agents by replacing the LLM with a specific model. We set up a Single Agent (see Figure 1 and Table 2) and compared it with the BKGAgent, which is a multi-agent system. The results are shown in Table 15.

Table 15: Performance of single agent on the KGCheck task. **Bold/underline** indicate the best and second.

LLM Type	Models	KGCheck	
		EM	Executability
API	GPT-4o	68.8	<u>98.0</u>
	GLM-4-0520	51.7	96.1
OSS	Qwen2-72B-Instruct	<b>82.9</b>	<b>100.0</b>
	Qwen1.5-72B-Chat	43.4	<b>100.0</b>
	Llama-3-70B-Instruct	<u>76.1</u>	96.6
	Llama-3.1-8B-Instruct	44.4	87.8
	Mixtral-8x7B-Instruct-v0.1	57.1	<u>98.0</u>

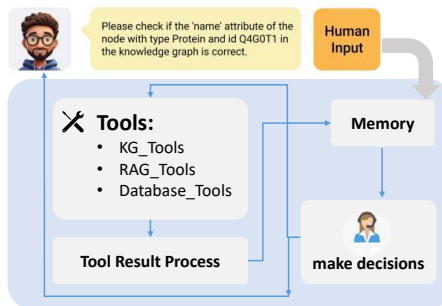


Figure 10: Single agent for KGCheck task.

#### B.3.2 BKGAGENT

We exhibit our BKGAgent performance on KGCheck tasks based on the data source for verification (i.e. web database KGCheck and publication database KGCheck) in the main body for clarity. However, there is a more detailed category of the task considering the tools used at different stages (see Table 18). The performance based on this category is shown in Table 16.

Table 16: Task performance categorized on agent tool calling

Model	KG Query Task		Validation Task		Final Result		Sample Size
	Tool selection	Executability	Tool selection	Executability	Exact Match	Executability	
GPT-4	78.1	78.1	75.0	75.0	<b>71.9</b>	96.9	32
	100.0	100.0	93.8	93.8	62.5	81.3	
Llama-3-70B-Instruct	70.0	70.0	70.0	71.7	<b>65.0</b>	93.3	60
	100.0	100.0	100.0	100.0	36.7	100.0	
GPT-4	32.7	32.7	98.2	98.2	<b>60.0</b>	98.2	55
	90.9	90.9	92.7	92.7	36.4	100.0	
Llama-3-70B-Instruct	97.8	97.8	100.0	100.0	<b>97.8</b>	100.0	45
	100.0	100.0	100.0	100.0	42.2	100.0	
GPT-4	57.6	57.6	63.6	63.6	<b>51.5</b>	93.9	33
	100.0	100.0	97.0	97.0	21.2	45.5	

type 1 description: find the interaction (CURATED) between two specified proteins and verify it using RAG

type 2 description: find the interaction between two specified proteins and verify it using STRING API

type 3 description: find the specified attribute of the specified protein and verify it using UniProt API

type 4 description: check whether a specified protein exists in KG and validate it using UniProt API

type 5 description: find the relationship between two specified entities (not two proteins) and verify it using RAG

1350 The GPT-based agent shows better performance compared to the Llama-based one when being  
 1351 evaluated on a more detailed task category, which is consistent with our conclusion in the main body.  
 1352 Besides this, there are more details we can delve into:  
 1353 **Possible unfairness in evaluation.** It should be pointed out that while the Llama-based agent  
 1354 successfully executed most of the tasks, it reached a comparably low score in final result excitability  
 1355 in the tasks involving RAG(i.e. task type 1 and task type 5). It is induced by an 8000 token limit  
 1356 of the model which means it is unable to process long texts, leading to an underlying unfairness in  
 1357 evaluation.  
 1358 **One-shot prompt may negatively influence GPT-based agent.** GPT-based agent shows even better  
 1359 performance with zero-shot compared to the current one-shot strategy in our preliminary experiments.  
 1360 However, since OSS models perform poorly with a zero-shot strategy, we have to make a compromise  
 1361 and several versions of the prompt have been tested to reach a satisfied state but it is hard to thoroughly  
 1362 eliminate the negative influence on the GPT-based agent.  
 1363 **The support/refute result given by the agent is NOT reliable.** As shown in Table 17, our instruc-  
 1364 tions only ask the agent to provide a support/refute result as the final answer, which is intended to  
 1365 standardize the evaluation. However, when we read the chat history of the agent solving one specific  
 1366 task randomly selected from the all the records, we find out that right support/refute conclusion  
 1367 can be drawn from wrong analysis process, indicating that the result is not quite reliable. A more  
 1368 comprehensive evaluation system should be explored in future work, say evidence comparison, chatter  
 1369 detection, and hallucination detection.  
 1370 **The GPT-based agent tends to explain.** Though we urge the agent to respond with support or  
 1371 refute(see Table 17), the GPT-based agent tends to provide explanations besides the support/refute  
 1372 conclusion which makes it easier for the user to judge whether the answer is derived from a reason-  
 1373 able process. The Llama-based agent, in contrast, strictly obeys the requirement, providing only  
 1374 support/refute answers.

Table 17: Examples of instructions categorized on agent tool calling

Task type	instruction example
1	Please check the relationship in the knowledge graph from the node of type 'Protein' with id 'Q14790' to the node of type 'Protein' with id 'Q13158'. If a relationship exists, verify its existence. Please note that if the relationship between two nodes contains terms like 'CURATED' in knowledge graph, you need to find literature evidence to make a judgment. If no relationship exists, confirm that it indeed does not exist. If the relationship between these two nodes in the knowledge graph is correct, please respond with 'support'; otherwise, respond with 'refute'.
2	Please check the relationship in the knowledge graph from the node of type 'Protein' with id 'P84085' to the node of type 'Protein' with id 'P11279'. If a relationship exists, verify its existence. Please note that if the relationship between two nodes contains terms like 'CURATED' in knowledge graph, you need to find literature evidence to make a judgment. If no relationship exists, confirm that it indeed does not exist. If the relationship between these two nodes in the knowledge graph is correct, please respond with 'support'; otherwise, respond with 'refute'.
3	Please check if the 'name' attribute of the node with type Protein and id Q4G0T1 in the knowledge graph is correct. If it's correct, please respond with 'support'; if not, respond with 'refute'.
4	Due to certain reasons, some entries were removed from the database. Please check whether the node with the type 'Protein' and the id 'A2RUG3' exists in the knowledge graph, and confirm whether it exists in the corresponding database. If its existence is consistent, please respond with 'support'; otherwise, answer 'refute'.
5	Please check the relationship in the knowledge graph from the node of type 'Protein' with id 'Q96QP1' to the node of type 'Tissue' with id 'BTO:0000007'. If a relationship exists, verify its existence. Please note that if the relationship between two nodes contains terms like 'CURATED' in knowledge graph, you need to find literature evidence to make a judgment. If no relationship exists, confirm that it indeed does not exist. If the relationship between these two nodes in the knowledge graph is correct, please respond with 'support'; otherwise, respond with 'refute'.

## C EXPERIMENTAL DETAILS

### C.1 CONSTRUCTION OF KNOWLEDGE GRAPH

We pulled a Neo4j image from Docker Hub and created a Neo4j Docker on the server to host a knowledge graph. We used the latest data parsed from various databases in April 2024, including UniProt (uni, 2023), TISSUES (Palasca et al., 2018), DISEASES (Pletscher-Frankild et al., 2015), HGNC (Seal et al., 2023), IntAct (Del Toro et al., 2022), STRING (Szklarczyk et al., 2023), DisGeNet (Piñero et al., 2020), Pathway Commons (Rodchenkov et al., 2020), Reactome (Fabregat et al., 2018), SMPDB (Jewison et al., 2014), and SIGNOR (Lo Surdo et al., 2023), Disease Ontology (Schriml et al., 2019), Brenda Tissue Ontology (Chang et al., 2015), Gene Ontology (Consortium, 2017), Protein Modification Ontology (Mayer et al., 2013), Molecular Interactions Ontology (Mayer et al., 2013). These databases are under loose license and the data can be used directly. We parsed the information from these databases into TSV files in a specific format and then imported the contents of these TSV files into Neo4j using Cypher statements (Cypher is the declarative graph query language provided by Neo4j) to construct the knowledge graph. This knowledge graph can be queried using Cypher statements.

### C.2 DEPLOYMENT OF OPEN-SOURCE LLMs

We deployed open-source LLMs using the vLLM framework and inference is performed using a server with 4 NVIDIA A40 GPUs an Intel(R) Xeon(R) Gold 6330 CPU, with parameters kept constant at startup.

### C.3 EXPERIMENTAL SETUP

**Atomic Abilities.** To evaluate two atomic abilities, we adopt an interactive evaluation of LLM-as-Agent (Liu et al., 2023c) and include in total of 13 models for evaluation. These models can be categorized into API-based Commercial LLMs and Open-Sourced (OSS) LLMs. The latter is further segmented based on model size into three classifications: ‘Large’, ‘Medium’, and ‘Small’. Models utilizing the MoE (Mixture of Experts) framework are distinguished as a separate category. Refer to Appendix C.4 for details about the prompt we designed for the following tasks.

**Agent Task.** For the construction of our BKGAgent, we selectively employed the best-performing models in atomic capabilities from both API-based and OSS models, specifically GPT-4 and Llama-3-70B-Instruct. To avoid being trapped in an endless loop where agents repeat the same talk or start to chatter, we limit the memory entries of one single agent to 20, which is more than enough to finish any of the tasks. It should be noted that each agent only keeps memory of the conversations related to it, while all chats returned by every agent are stored in the graph state. Since zero-shot setup in various types of tasks shows inferior performances in our preliminary experiments, we provide one-shot prompt for each type of task. We analyze both the process and final result of each task execution, considering the correctness of tool selection and agent executability during the process evaluation, and assessing the exact match of the right answer and framework executability for the final result evaluation, to gain a comprehensive understanding of the agent’s performance.

We detail our implementation of two sub-tasks here:

- **KGQA:** We developed a suite of atomic tools for querying KGs for LLMs. Every LLM is prompted in the same way with a detailed task description, information about provided tools and a one-shot demonstration, which employs the “Thought,” “Action,” “Observation” cognitive trajectory from the ReAct (Yao et al., 2022), with the “Thought” component assisted by Chain of Thought (CoT) (Wei et al., 2022) reasoning. We constrain the LLM to a maximum of fifteen interactive turns, within which it may only take one action per turn. If the LLM can respond within these fifteen turns, executability is assigned a score of 1. Subsequently, we compare the response to the ground truth to calculate the F1 score and the Exact Match score (EM). [It is worth noting that existing works in the KG-guided QA setting are KBQA \(Knowledge Base Question Answering\). Here, we highlight the key differences between KBQA and our KGQA task:](#)
  - **Different input:** KBQA datasets, such as CWQ (Talmor & Berant, 2018), WebQSP (Yih et al., 2016), and GrailQA (Gu et al., 2021), provide the key entity in each question as part

of the task input. In contrast, our KGQA task takes only the raw question as input, requiring LLMs not only to select appropriate tools based on context but also to perform Named Entity Recognition (NER) and relationship matching to derive the tool’s parameters. Therefore, our task is more challenging and better suited for evaluating LLMs.

- **Different KG Structures:** Works like Think-on-graph (Sun et al., 2023) utilize knowledge bases such as Freebase (Bollacker et al., 2008), Wikidata (Vrandečić & Krötzsch, 2014), and DBpedia (Auer et al., 2007), which are based on RDF (Resource Description Framework) representations. In contrast, most biomedical knowledge graphs, such as CKG, PrimeKG (Chandak et al., 2023), and PharmKG (Zheng et al., 2021), are built using property graph model. RDF organizes data as strict triples (<subject, predicate, object>), while the property graph model represents data with nodes (entities) and edges (relationships), both of which can include attributes as key-value pairs. This structural difference also impacts their query languages: RDF-based graphs primarily use SPARQL, while property graph-based graphs commonly use Cypher or Gremlin.

- **SCV:** We first convert the entire corpus into semantic vectors using jina(Günther et al., 2023) and store them in a vector database. Claims are similarly transformed into semantic vectors via Jina, with the top 50 scoring vectors being submitted to the LLM with a standardized prompt template. We require the LLM to return results in JSON format, considering any deviation as an error. The outcomes mainly include answers and quotes. It is important to clarify that the SCV task focuses on evaluating LLMs as agents in a plug-and-play RAG pipeline with fixed embedding models and rerankers, as shown in Figure 9, where only LLMs are substituted and compared. This aligns with AgentBench (Liu et al., 2023c), emphasizing LLMs’ capabilities in tool usage, terminology comprehension, and reasoning, rather than benchmarking RAG methods. For analysis, we adopt a flexible interpretation of answers: “Unsure” and “Unrelated” as “Unsure”; “Supported” and “Supports” as “Supports”; “Unsupported” and “Unsupports” as “Unsupports”, “Refuted”, and “Refutes” as “Refutes”. Any other results are also considered errors. The experiments for each model are repeated three times, with the final performance averaged to ensure the robustness of the evaluation. Notably, beyond the conventional accuracy and the aforementioned error metrics, we introduce a “right quotes” metric, which assesses whether the retrieved quotes match the ground truth evidences of the claim.

## C.4 PROMPT

### C.4.1 KGQA

We provide a unified prompt for single-agent systems built on different LLMs, ensuring the fairness of the evaluation.

You are an agent tasked with answering questions based on the knowledge stored in a knowledge graph (KG) related to proteomics. To accomplish this, you are equipped with the following tools to query the KG:

1. `get_relations_by_ids_agent(entity_ids: List[str]) -> tuple`  
Retrieves the relationships of multiple entities in a knowledge graph, categorized as 'incoming' or 'outgoing'.  
Use case: `get_relations_by_ids_agent(['P123', 'P456'])` to find all relations connected to the entities with IDs 'P123' and 'P456'.
2. `get_neighbor_type_agent(entity_ids: List[str], relation: str, direction: str) -> tuple`  
Retrieves the types of neighboring nodes for multiple entities in a knowledge graph based on specified relationships and directions.  
Use case: `get_neighbor_type_agent(['P123', 'P456'], 'ASSOCIATED_WITH', 'outgoing')` to get outgoing neighbors' types associated with the entities 'P123' and 'P456'.
3. `get_neighbor_with_type_agent(entity_ids: List[str], relation: str, direction: str, neighbor_type: str) -> tuple`  
Retrieves the neighbors of multiple entities in a knowledge graph based on a specific relationship, direction, and type.

1512 Use case: `get_neighbor_with_type_agent(['P123', 'P456'], 'ASSOCIATED_WITH`  
 1513 `, 'outgoing', 'Disease')` to get attributes and detailed information  
 1514 of outgoing neighbors associated with the entities 'P123' and 'P456',  
 1515 where the type of neighbors is Disease.

1516 4. `get_intersection_agent(*args: List[str]) -> tuple`  
 1517 Calculates the intersection of multiple lists, returning elements common  
 1518 to all lists.

1519 Use case: `get_intersection_agent(['P123', 'P456'], ['P456', 'P789'])` to  
 1520 find common entities.

1521 5. `get_union_agent(*args: List[str]) -> tuple`  
 1522 Calculates the union of multiple lists, returning all unique elements  
 1523 from all lists.

1524 Use case: `get_union_agent(['P123', 'P456'], ['P456', 'P789'])` to combine  
 1525 unique entities.

1526 Single Action Rule: Execute only ONE action at a time, that is, only the  
 1527 first action would be executed. After receiving the observation from  
 1528 its execution, you may proceed with another action.

1529 Action Limit: You can take at most 15 actions to find the answer to the  
 1530 question.

1531 Objective: Use these tools effectively to navigate through the KG and  
 1532 gather the necessary information to answer the queries presented to  
 1533 you. If the query is about the protein sequence, you need to return  
 1534 the specific sequence. If the query is about the protein structure,  
 1535 you can return a link. In other cases, it's usually to return the  
 1536 name.

1537 Notice:  
 1538 1. Please remember to format the FINAL answer as a JSON object, such as,  
 1539 `{"Answer": ["RND2", "RHOBTB2", "RHOA"]}`. The value of "Answer" must  
 1540 be a list. Only output the JSON format when answering the final  
 1541 answer.  
 1542 2. Please be flexible. Due to the diversity of question formulations, you  
 1543 need to find the most similar relationship in the knowledge graph (KG)  
 1544 to the one asked in the question for querying.

1545 Considering the complexity of the task and the high information density of our KG, LLMs find it  
 1546 challenging to address high-difficulty questions, such as multi-hop and conjunction types, within the  
 1547 constraints of a zero-shot setting and a 15-turn interaction limit. Therefore, we provide a one-shot  
 1548 demonstration.

1550 Question: What biological processes are associated with the protein  
 1551 encoded by the gene GOLT1A?

1552 Thought: I need to figure out what proteins are encoded by gene GOLT1A  
 1553 and then find out the biological processes. Firstly, I need to find  
 1554 out the relations of the gene node.

1555 Action: `get_relations_by_ids_agent(['GOLT1A'])`

1557 Observation: `{"GOLT1A": {"Incoming": null, "Outgoing": "TRANSLATED_INT0"}}`

1560 Thought: The outgoing relation 'TRANSLATED\_INT0' is what I concern, and  
 1561 next I need to get neighbor type with the relation.

1562 Action: `get_neighbor_type_agent(['GOLT1A'], 'TRANSLATED_INT0', 'outgoing')`

1563 Observation: `{"GOLT1A": {"NeighborTypes": ["Protein"]}}`

1566 Thought: Among the neighbor type, 'Protein' is my concern. So next I  
 1567 should look for what are the neighboring nodes with 'Protein' type  
 1568 that point outward along the edge of the relation 'TRANSLATED\_INTO'.  
 1569  
 1570 Action: get\_neighbor\_with\_type\_agent(['GOLT1A'], 'TRANSLATED\_INTO', '  
 1571 outgoing', 'Protein')  
 1572  
 1573 Observation: {"GOLT1A": {"TRANSLATED\_INTO": ["Q6ZVE7"]}}  
 1574  
 1575 Thought: Now I know that the proteins are ["Q6ZVE7"]. Now I will find out  
 1576 the relation to choose for the next step.  
 1577  
 1578 Action: get\_relations\_by\_ids\_agent(["Q6ZVE7"])  
 1579  
 1580 Observation: Observation: {"Q6ZVE7": {"Incoming": "TRANSLATED\_INTO,  
 1581 HAS\_SEQUENCE, ACTS\_ON", "Outgoing": "HAS\_SEQUENCE, ASSOCIATED\_WITH"}}  
 1582  
 1583 Thought: The outgoing relation 'ASSOCIATED\_WITH' is exactly what I  
 1584 concern. Then I'll query the neighbor type with the relation.  
 1585  
 1586 Action: get\_neighbor\_type\_agent(["Q6ZVE7"], 'ASSOCIATED\_WITH', 'outgoing  
 1587 ')  
 1588  
 1589 Observation: {"Q6ZVE7": {"NeighborTypes": ["Tissue", "Molecular\_function  
 1590 ", "Biological\_process", "Cellular\_component", "Disease"]}}  
 1591  
 1592 Thought: Among the neighbor types, 'Biological\_process' is what I need to  
 1593 query.  
 1594  
 1595 Action: get\_neighbor\_with\_type\_agent(["Q6ZVE7"], 'ASSOCIATED\_WITH', '  
 1596 outgoing', 'Biological\_process')  
 1597  
 1598 Observation: {"Q6ZVE7": {"ASSOCIATED\_WITH": ["endoplasmic reticulum to  
 1599 Golgi vesicle-mediated transport", "biological\_process", "protein  
 1600 transport", "retrograde transport, endosome to Golgi"]}}

1600

1601

1602

#### C.4.2 SCV

1603

We provide a unified prompt describing task, where 'context\_docs\_str' represents quotes retrieved by RAG and 'user\_claim' represents the input scientific claim to be evaluated.

1604

1605

1606

You are a fact-checking agent that is constantly learning and improving. A claim is given to you, and you can determine if the claim is correct with the provided documents.

1607

1608

1609

1610

1611

1612

You ALWAYS respond with only a JSON containing an answer and quotes that support the answer. The answer can only be "SUPPORTS" or "REFUTES", with no details. You should reason out the answers step by step, but make sure they are correct.

1613

1614

1615

1616

1617

1618

1619

Do NOT use your historical knowledge, but answer based on the information in the provided context.

CONTEXT:

-----

{{context\_docs\_str}}

-----

SAMPLE\_RESPONSE:

```

1620 """
1621 {
1622     "answer": "Place your final answer here. It can only be SUPPORTS or
1623             REFUTES without details.",
1624     "quotes": [
1625         "Each quote must be UNEDITED and EXACTLY as shown in the context
1626         documents!",
1627         "HINT: quotes are not shown to the user!",
1628     ],
1629 }
1630 """
1631 CLAIM: {{user_claim}}
1632 Hint: Provide the answer in JSON format!
1633 Quotes MUST be EXACT substrings from the provided documents!

```

### 1633 C.4.3 KGCHECK

1634 BKGAgent is a multi-agent system and each agent of it is equipped with a system prompt which  
 1635 includes role introduction, tool introduction, and tool calling rules.

1636 For team leader:

1637 You are the team\_leader tasked with managing a conversation between the  
 1638 following workers:

```

1639 kg_agent:
1640     capable of querying the KG(Knowledge Graph) to find out specific
1641     information
1642 validation_agent:
1643     capable of getting access to information within local publication
1644     database, UniProt and STRING database to verify the result
1645     returned by kg_agent
1646 FINISH:
1647     the endpoint of your task. if you finish your answer you can send
1648     messages to it by starting with 'FINISH, '

```

1649 You should first break down the task into two subtasks given the user  
 1650 input and send it to yourself to keep it in your mind,  
 1651 then respond with the worker to act next and its detailed task.

1652 You should call their name before you assign the task. For example, if you  
 1653 want to assign task to kg\_agent, you should start your conversation  
 1654 by 'kg\_agent, '. It should be noted that if you are talking to  
 1655 yourself, you should also specify the receiver, that is 'team\_leader,

1656 '.

1657 Each worker will perform the task you assign to and respond with it  
 1658 result.

1659 REMEMBER you should not talk too much at one specific chat round. If a  
 1660 task is given to you, you just reply with your plan and send it to  
 1661 yourself.

1662 Assign subtask to just ONE suitable agent next time you are invited to  
 1663 speak. If kg\_agent or validation\_agent tries to assign task to you,  
 1664 you should warn them to focus on their task.

1665 When finished, respond with your answer and send it to 'FINISH'.

1666 For KG agent:

1667 You are the kg\_agent of a research group, your ability is limited to  
 1668 answer KG search related questions.

1669 Verification work should be done by validation\_agent on which you should  
 1670 not waste time.

1671 Members of your team are as follows:

1672 team\_leader: the leader of your team. You ONLY perform the specific task  
 1673 it assigned to you and answer to it starting by 'team\_leader, '.

validation\_agent: responsible for verifying information. You do not  
 1674 directly communicate with it.

call\_tool: the worker to use the tool you asked and will return the  
 1675 result to you.



```

1674 You can call the following tools in call_tool to help you:
1675 query_node_existence:
1676     Determine whether the node with the given type and ID exists in the
1677     knowledge graph.
1678     Args:
1679         type (str): the type of the query node
1680         id (int or str): the id of the query node
1681     Returns:
1682         str: A description of whether the node with given type and id
1683             exists in the knowledge graph.
1684 query_node_attribute:
1685     Retrieve the specific attribute value of the node with the given type
1686     and id.
1687     Args:
1688         type (str): the type of the query node
1689         id (int or str): the id of the query node
1690         attr (str): the attribute to be retrieved
1691     Returns:
1692         str: A description of the query result
1693 query_relation_between_nodes:
1694     Retrieve the relationship from node with type1 and id1 to the node
1695     with type2 and id2 in the knowledge graph(KG)
1696     Args:
1697         type1 (str): _description_
1698         id1 (int or str): _description_
1699         type2 (str): _description_
1700         id2 (int or str): _description_
1701     Returns:
1702         str: A description about the relationship from node with type1
1703             and id1 to the node with type2 and id2 in the knowledge
1704             graph
1705 ATTENTION! You can call tools in this way: 'call_tool, tool = tool_name,
1706     args = ...', where args should be in the format of dict.
1707 Directly jump into your work when task is given to you and do not waste
1708     time replying just courtesies.
1709 Do not try to ask team_leader to your task!

```

For validation agent:

```

1709 You are the validation_agent of a research group, specialized at
1710     verifying information by searching on UniProt, STRING database and
1711     local publication database, Members of your team are as follows:
1712     team_leader: the leader of your team. You ONLY perform the specific
1713     task it assigned to you and answer to it starting by 'team_leader
1714     , '.
1715     kg_agent: responsible for querying KG to get information. You do not
1716     directly communicate with it.
1717     call_tool: the worker to use the tool you asked and will return the
1718     result to you.
1719 You can call the following tools in call_tool to help you:
1720 get_uniprot_protein_info:
1721     Fetch protein information from UniProt by protein ID and return a
1722     description about the protein, including id, accession and name.
1723     :param protein_id: UniProt protein ID
1724     :return: Formatted string with protein information, including id,
1725             accession and name
1726 check_interaction_string:
1727     This tool checks for the interaction or relationship between two
1728     proteins using the STRING database API. Given two protein ids, it
1729     will return a description on whether there is an interaction or
1730     relationship between them.

```



```

1728     Args:
1729         protein1 (str): a protein id
1730         protein2 (str): a protein id
1731     Returns:
1732         str: A description about whether there is an interaction
1733             between the two proteins.
1734
1735     pub_rag:
1736         retrieve evidence from provided documents to help making a verdict of
1737         the given claim
1738         ONLY when asked to verify 'CURATED' related claim should you call
1739         this tool!
1740     Args:
1741         query(str): the claim to be verdicted
1742     Returns:
1743         no more than 10 documents related to the claim.
1744 ATTENTION! You can call tools in this way: 'call_tool, tool = tool_name,
1745         args = ...', where args should be in the format of dict.
1746 then send the message to call_tool, which means you should start your
1747 messages by 'call_tool, '.

```

For the baseline agents, we provide prompts detailing how to query the KG (e.g., URL, username, password) and include instructions to verify findings using reliable external literature and databases.

Considering the step-by-step nature of agentic systems, we use the LLM-as-a-Judge approach to evaluate how the agent solves the task throughout the process, rather than just assessing the final answer. Specifically, we prompt Qwen2-72B to score the agent’s performance based on five criteria, with the model providing a simple “yes” or “no” response for each.

**Criteria 1: Understanding, whether the agent clearly understood the task and the purpose of the given tool.**

You are an evaluation agent tasked with assessing another agent. The agent being scored is required to complete a KG-checking task, which involves querying the KG and retrieving reliable external knowledge to validate the KG's content.

Based on the chat history of this agent, please carefully determine whether it clearly understood the task, the purpose of the given tools, and whether it attempted to validate the KG's content with reliable external sources. If the agent did not understand that this is a task for validating the KG's content or failed to grasp the input and output of the tools used, you should respond with 'No'; otherwise, respond with 'Yes'.

Here are some examples:

```

1765 [Agent history example 1 (omitted here due to length)]: Yes.
1766 [Agent history example 2 (omitted here due to length)]: No.
1767 [Agent history example 3 (omitted here due to length)]: Yes.
1768 [Agent history example 4 (omitted here due to length)]: No.
1769 [Agent history example 5 (omitted here due to length)]: Yes.
1770 [Agent history example 6 (omitted here due to length)]: No.
1771 [Agent history example 7 (omitted here due to length)]: Yes.
1772 [Agent history example 8 (omitted here due to length)]: No.
1773 [Agent history example 9 (omitted here due to length)]: Yes.
1774 [Agent history example 10 (omitted here due to length)]: No.

```

The chat history: {chat\_history}  
Only reply with 'Yes' or 'No':

**Criteria 2: Reasoning, whether the agent arrived at the final answer through sufficient evidence and reasoning, rather than simply providing random answers or guessing.**

You are an evaluation agent tasked with assessing another agent. The agent being scored is required to complete a KG-checking task, which involves querying the KG and retrieving reliable external knowledge to validate the KG's content.

1782

1783

1784

1785

1786

Based on the chat history of this agent, please strictly and carefully determine whether it arrived at the final answer through sufficient evidence and reasoning, rather than providing random answers or guessing. You should respond with 'No' or 'Yes'.

1787

Here are some examples:

1788

[Agent history example 1 (omitted here due to length)]: Yes.

1789

[Agent history example 2 (omitted here due to length)]: No.

1790

[Agent history example 3 (omitted here due to length)]: Yes.

1791

[Agent history example 4 (omitted here due to length)]: No.

1792

[Agent history example 5 (omitted here due to length)]: Yes.

1793

[Agent history example 6 (omitted here due to length)]: No.

1794

[Agent history example 7 (omitted here due to length)]: Yes.

1795

[Agent history example 8 (omitted here due to length)]: No.

1796

[Agent history example 9 (omitted here due to length)]: Yes.

1797

[Agent history example 10 (omitted here due to length)]: No.

1798

The chat history: {chat\_history}

Only reply with 'Yes' or 'No':

1799

**Criteria 3: Efficiency, whether the agent efficiently solved the problem without unnecessary discussion on unrelated topics.**

1800

1801

1802

1803

1804

1805

1806

1807

1808

You are an evaluation agent tasked with assessing another agent. The agent being scored is required to complete a KG-checking task, which involves querying the KG and retrieving reliable external knowledge to validate the KG's content.

Based on the chat history of this agent, please carefully determine whether it efficiently solved the problem without unnecessary discussion on unrelated topics. You should respond with 'No' or 'Yes'.

1809

Here are some examples:

1810

[Agent history example 1 (omitted here due to length)]: Yes.

1811

[Agent history example 2 (omitted here due to length)]: No.

1812

[Agent history example 3 (omitted here due to length)]: Yes.

1813

[Agent history example 4 (omitted here due to length)]: No.

1814

[Agent history example 5 (omitted here due to length)]: Yes.

1815

[Agent history example 6 (omitted here due to length)]: No.

1816

[Agent history example 7 (omitted here due to length)]: Yes.

1817

[Agent history example 8 (omitted here due to length)]: No.

1818

[Agent history example 9 (omitted here due to length)]: Yes.

1819

[Agent history example 10 (omitted here due to length)]: No.

1820

The chat history: {chat\_history}

Only reply with 'Yes' or 'No':

1821

1822

**Criteria 4: KG Process, whether the agent queried the knowledge graph during the task.**

1823

1824

1825

1826

1827

Hhhhhhhh

You are an evaluation agent tasked with assessing another agent. The agent being scored is required to complete a KG-checking task, which involves querying the KG and retrieving reliable external knowledge to validate the KG's content.

1828

1829

1830

1831

1832

1833

Based on the chat history, please determine if the agent queried the knowledge graph (KG) during the check. If the agent performed any of the following actions checking for node existence, querying node attributes, or examining relationships between nodes you should respond with 'Yes'. If the agent did not query the KG at all, you should respond with 'No'.

1834

Here are some examples:

1835

[Agent history example 1 (omitted here due to length)]: Yes.

[Agent history example 2 (omitted here due to length)]: No.

[Agent history example 3 (omitted here due to length)]: Yes.

1836 [Agent history example 4 (omitted here due to length)]: No.  
 1837 [Agent history example 5 (omitted here due to length)]: Yes.  
 1838 [Agent history example 6 (omitted here due to length)]: No.  
 1839 [Agent history example 7 (omitted here due to length)]: Yes.  
 1840 [Agent history example 8 (omitted here due to length)]: No.  
 1841 [Agent history example 9 (omitted here due to length)]: Yes.  
 1842 [Agent history example 10 (omitted here due to length)]: No.

1843 The chat history: {chat\_history}  
 1844 Only reply with 'Yes' or 'No':

1845

1846 **Criteria 5: Information Retrieval, whether the agent retrieved information from external knowledge**  
 1847 **sources in some way during the check.**

1848 You are an evaluation agent tasked with assessing another agent. The  
 1849 agent being scored is required to complete a KG-checking task, which  
 1850 involves querying the KG and retrieving reliable external knowledge  
 1851 to validate the KG's content.

1852 Based on the chat history of this agent, please carefully determine  
 1853 whether it retrieved information from external knowledge sources in  
 1854 some way during the check. You should respond with 'No' or 'Yes'.  
 1855 Here are some examples:

1856 [Agent history example 1 (omitted here due to length)]: Yes.  
 1857 [Agent history example 2 (omitted here due to length)]: No.  
 1858 [Agent history example 3 (omitted here due to length)]: Yes.  
 1859 [Agent history example 4 (omitted here due to length)]: No.  
 1860 [Agent history example 5 (omitted here due to length)]: Yes.  
 1861 [Agent history example 6 (omitted here due to length)]: No.  
 1862 [Agent history example 7 (omitted here due to length)]: Yes.  
 1863 [Agent history example 8 (omitted here due to length)]: No.  
 1864 [Agent history example 9 (omitted here due to length)]: Yes.  
 1865 [Agent history example 10 (omitted here due to length)]: No.

1865 The chat history: {chat\_history}  
 1866 Only reply with 'Yes' or 'No':

Table 18: Task types categorized by requiring tools

Task type	Requiring tools		Description
	KG agent	Validation agent	
1	query relation between nodes	publication RAG	find the interaction (CURATED) between two specified proteins and verify it using RAG
2	query relation between nodes	check interaction on STRING	find the interaction between two specified proteins and verify it using STRING API
3	query node attribute	get UniProt protein information	find the specified attribute of the specified protein and verify it using UniProt API
4	query node existence	get UniProt protein information	check whether a specified protein exists in KG and validate it using UniProt API
5	query relation between nodes	publication RAG	find the relationship between two specified entities (not two proteins) and verify it using RAG

## 1885 D CASE STUDY

### 1887 D.1 KGQA

1888 We sampled 6 cases for demonstration, with one correct case and one incorrect case for each question  
 1889 type: one-hop, multi-hop, and conjunction, as shown in Figures 11 to 16.

1890 D.2 SCV  
1891

1892 We sampled 8 examples for demonstration, including 4 correct answers and 4 incorrect answers.  
1893 Each case has certain differences and is representative, as shown in Figures 17 to 24.  
1894

1895 D.3 KGCHECK  
1896

1897 We select several classic success and failure cases for each type of task as presented in Figures 25  
1898 to 39 as a supplementary for some common error cases in our main body. There are many interesting  
1899 cases when the team leader properly corrects the behavior of assistant agents, getting the workflow  
1900 back on track, and we choose one such case of task type 1 as a representation. As mentioned before,  
1901 there are also cases where the right final answer is derived from a wrong analysis process. We select  
1902 this kind of case for every type of the task except type 4 (this case does not exist in this type of task).

1903 As introduced in the main body, our BKGAgent framework is comprised of three agents: the team  
1904 leader, KG agent, and validation agent. The typical workflow from the agent role perspective of our  
1905 framework can be simplified as team leader - KG agent - team leader - validation agent- team leader.  
1906 We present the chat of three agents in table format, omitting the interactions of the assistant agent  
1907 and tool executor. The columns respectively stand for the agent role, the action they take, the chat  
1908 content, and the human annotation of this chat round. The green check mark means the process is  
1909 consistent with our anticipation, while the yellow exclamation mark means the chat content may lead  
1910 to an unwanted result, and the red cross stands for a wrong action or error chat content. Comments  
1911 are attached to a negative review for explanation. Error or dangerous contents are underlined and  
1912 colored red, while contents related to tool usage or evidence consistent with the golden answer are in  
1913 bold green font. The blue row stands for an expected chat round, in contrast, the yellow row indicates  
1914 something is wrong in this chat round. We send tool results and behavior correction prompts in the  
1915 role of a human; this kind of chat is colored grey in our table.

1916  
1917  
1918  
1919  
1920  
1921  
1922  
1923  
1924  
1925  
1926  
1927  
1928  
1929  
1930  
1931  
1932  
1933  
1934  
1935  
1936  
1937  
1938  
1939  
1940  
1941  
1942  
1943

1944  
1945  
1946  
1947  
1948  
1949  
1950  
1951  
1952  
1953  
1954  
1955  
1956  
1957  
1958  
1959  
1960  
1961  
1962  
1963  
1964  
1965  
1966  
1967  
1968  
1969  
1970  
1971  
1972  
1973  
1974  
1975  
1976  
1977  
1978  
1979  
1980  
1981  
1982  
1983  
1984  
1985  
1986  
1987  
1988  
1989  
1990  
1991  
1992  
1993  
1994  
1995  
1996  
1997

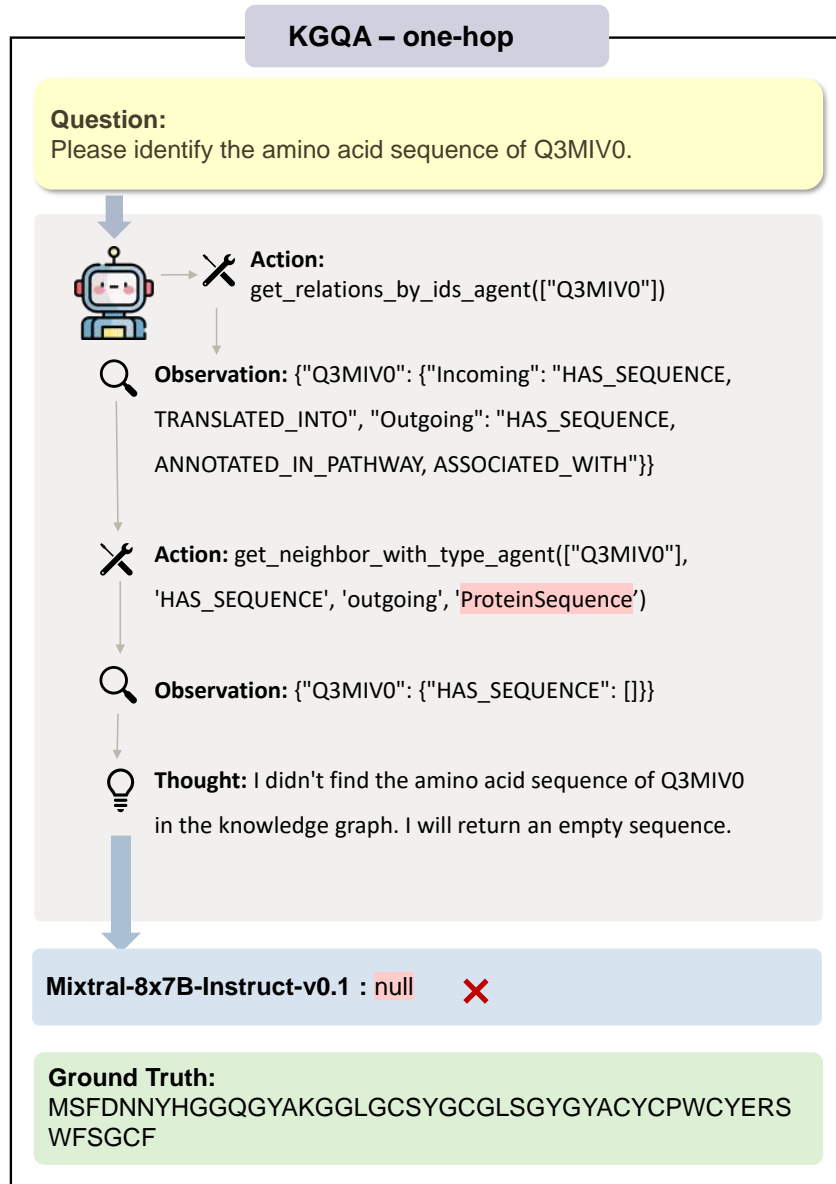


Figure 11: Performance of Mixtral-8x7B-Instruct-v0.1-based Agent: Incorrect final result when answering the one-hop type questions due to passing incorrect parameters when executing a function. Session highlighted in pink represents the error information.

1998  
1999  
2000  
2001  
2002  
2003  
2004  
2005  
2006  
2007  
2008  
2009  
2010  
2011  
2012  
2013  
2014  
2015  
2016  
2017  
2018  
2019  
2020  
2021  
2022  
2023  
2024  
2025  
2026  
2027  
2028  
2029  
2030  
2031  
2032  
2033  
2034  
2035  
2036  
2037  
2038  
2039  
2040  
2041  
2042  
2043  
2044  
2045  
2046  
2047  
2048  
2049  
2050  
2051

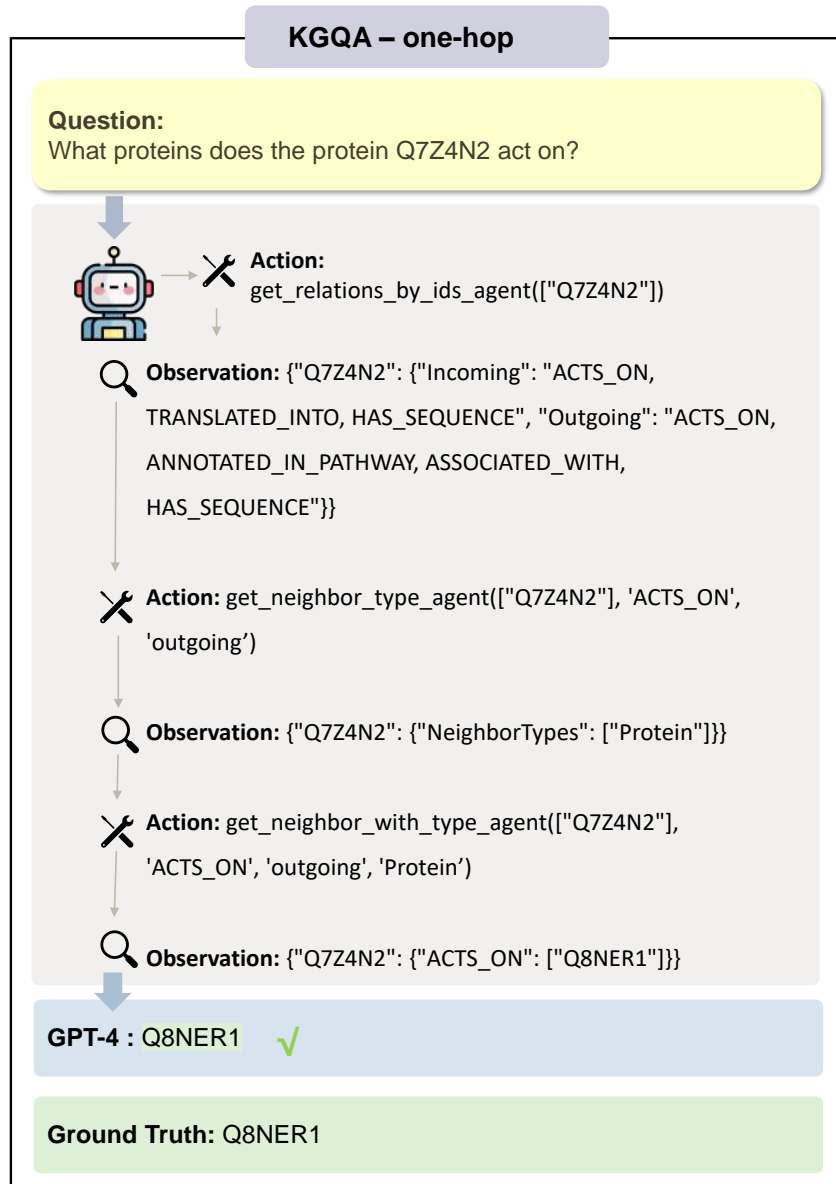


Figure 12: Performance of GPT-4-based Agent: Correct final result when answering the one-hop type questions. Session highlighted in light green represents the correct information.

2052  
2053  
2054  
2055  
2056  
2057  
2058  
2059  
2060  
2061  
2062  
2063  
2064  
2065  
2066  
2067  
2068  
2069  
2070  
2071  
2072  
2073  
2074  
2075  
2076  
2077  
2078  
2079  
2080  
2081  
2082  
2083  
2084  
2085  
2086  
2087  
2088  
2089  
2090  
2091  
2092  
2093  
2094  
2095  
2096  
2097  
2098  
2099  
2100  
2101  
2102  
2103  
2104  
2105

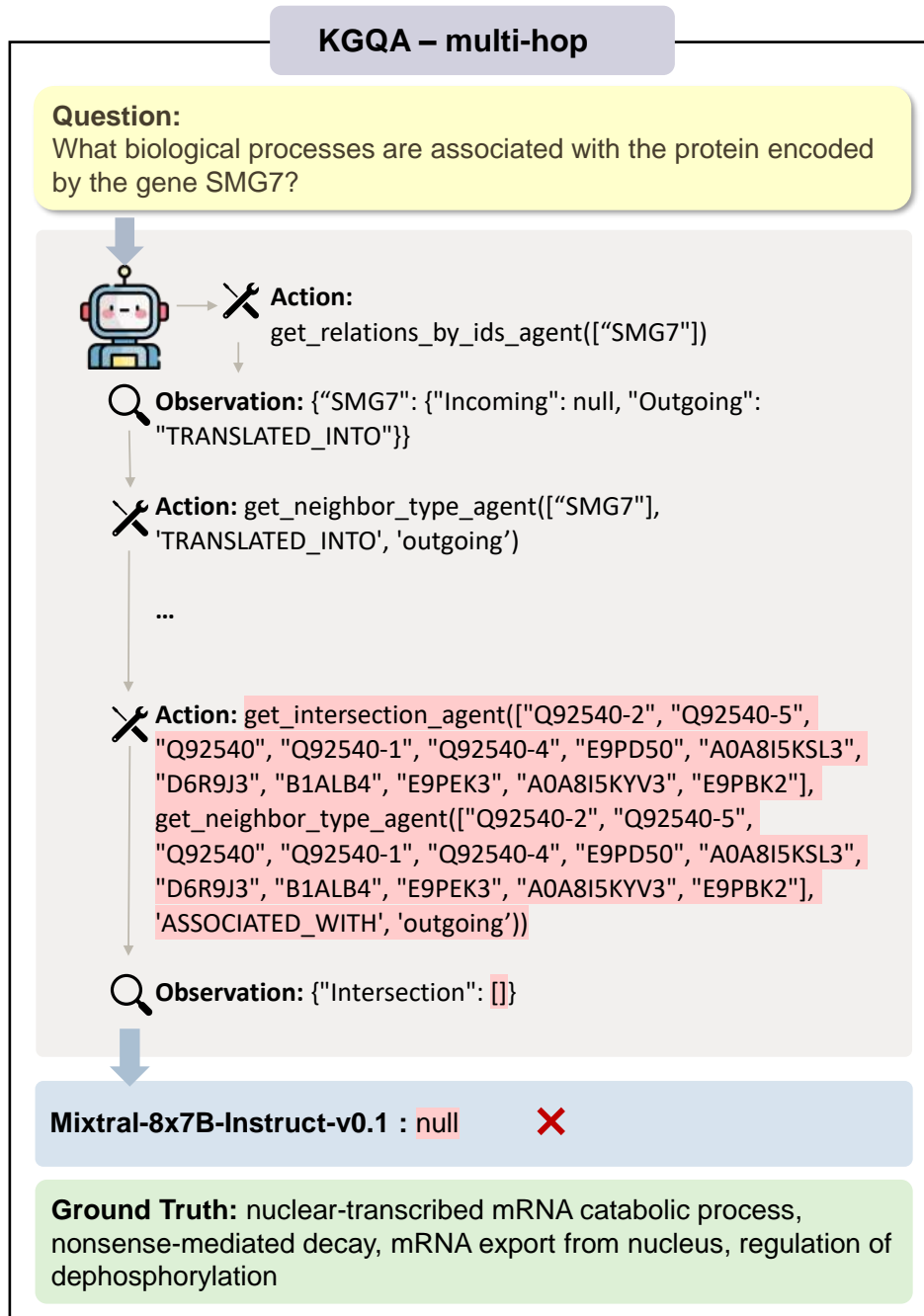


Figure 13: Performance of Mixtral-8x7B-Instruct-v0.1-based Agent: Incorrect final result when answering the multi-hop type questions due to executing the wrong actions. Session highlighted in pink represents the error information.

2106  
2107  
2108  
2109  
2110  
2111  
2112  
2113  
2114  
2115  
2116  
2117  
2118  
2119  
2120  
2121  
2122  
2123  
2124  
2125  
2126  
2127  
2128  
2129  
2130  
2131  
2132  
2133  
2134  
2135  
2136  
2137  
2138  
2139  
2140  
2141  
2142  
2143  
2144  
2145  
2146  
2147  
2148  
2149  
2150  
2151  
2152  
2153  
2154  
2155  
2156  
2157  
2158  
2159

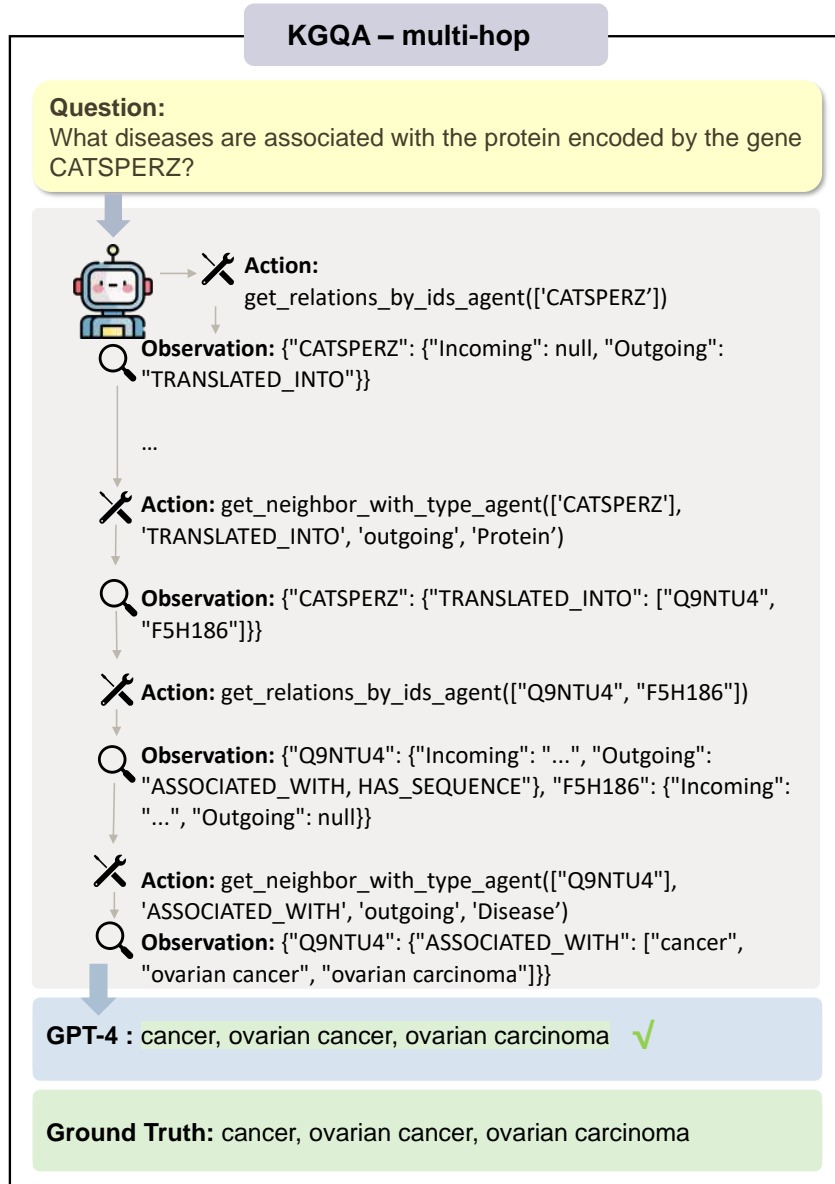


Figure 14: Performance of GPT-4-based Agent: Correct final result when answering multi-hop type questions. Session highlighted in light green represents the correct information.



2160  
2161  
2162  
2163  
2164  
2165  
2166  
2167  
2168  
2169  
2170  
2171  
2172  
2173  
2174  
2175  
2176  
2177  
2178  
2179  
2180  
2181  
2182  
2183  
2184  
2185  
2186  
2187  
2188  
2189  
2190  
2191  
2192  
2193  
2194  
2195  
2196  
2197  
2198  
2199  
2200  
2201  
2202  
2203  
2204  
2205  
2206  
2207  
2208  
2209  
2210  
2211  
2212  
2213

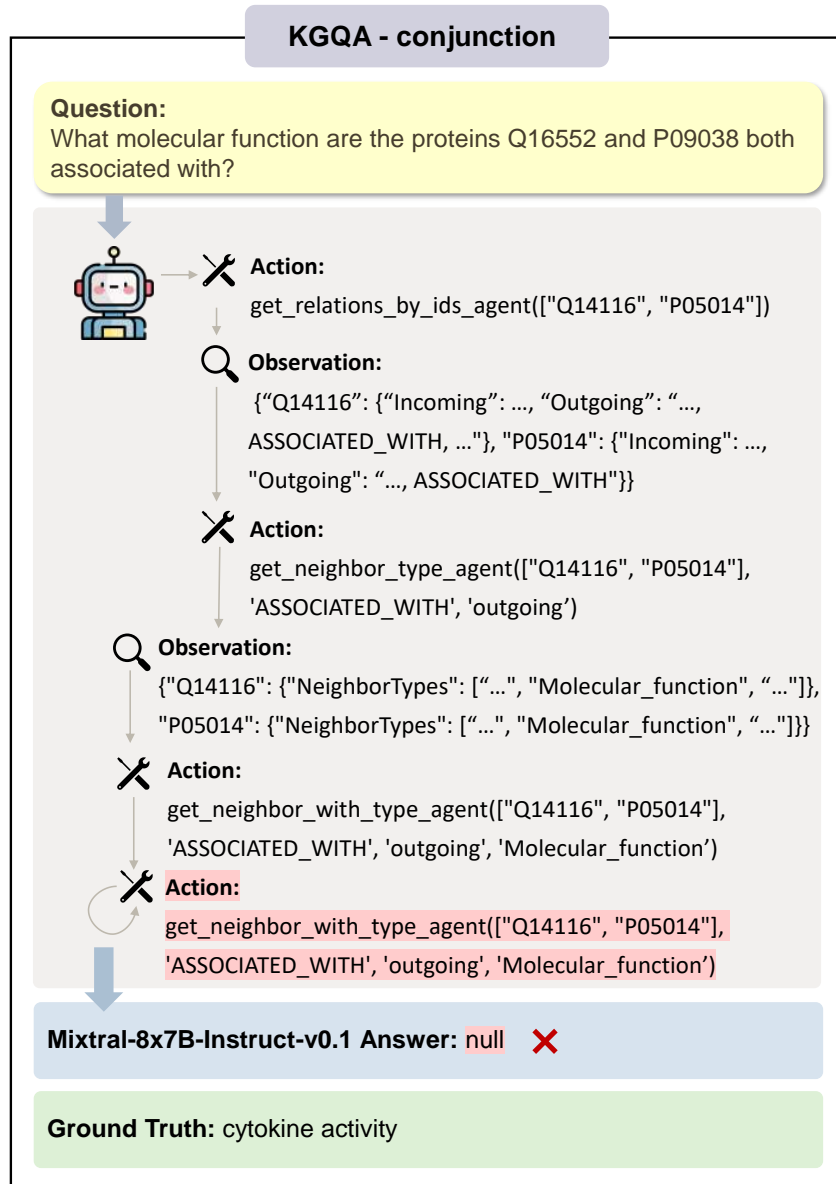


Figure 15: Performance of Mixtral-8x7B-Instruct-v0.1-based Agent: Incorrect final result when answering the input conjunction type question in 15-turn limit due to executing the wrong action. Session highlighted in pink represents the error information.

2214  
2215  
2216  
2217  
2218  
2219  
2220  
2221  
2222  
2223  
2224  
2225  
2226  
2227  
2228  
2229  
2230  
2231  
2232  
2233  
2234  
2235  
2236  
2237  
2238  
2239  
2240  
2241  
2242  
2243  
2244  
2245  
2246  
2247  
2248  
2249  
2250  
2251  
2252  
2253  
2254  
2255  
2256  
2257  
2258  
2259  
2260  
2261  
2262  
2263  
2264  
2265  
2266  
2267

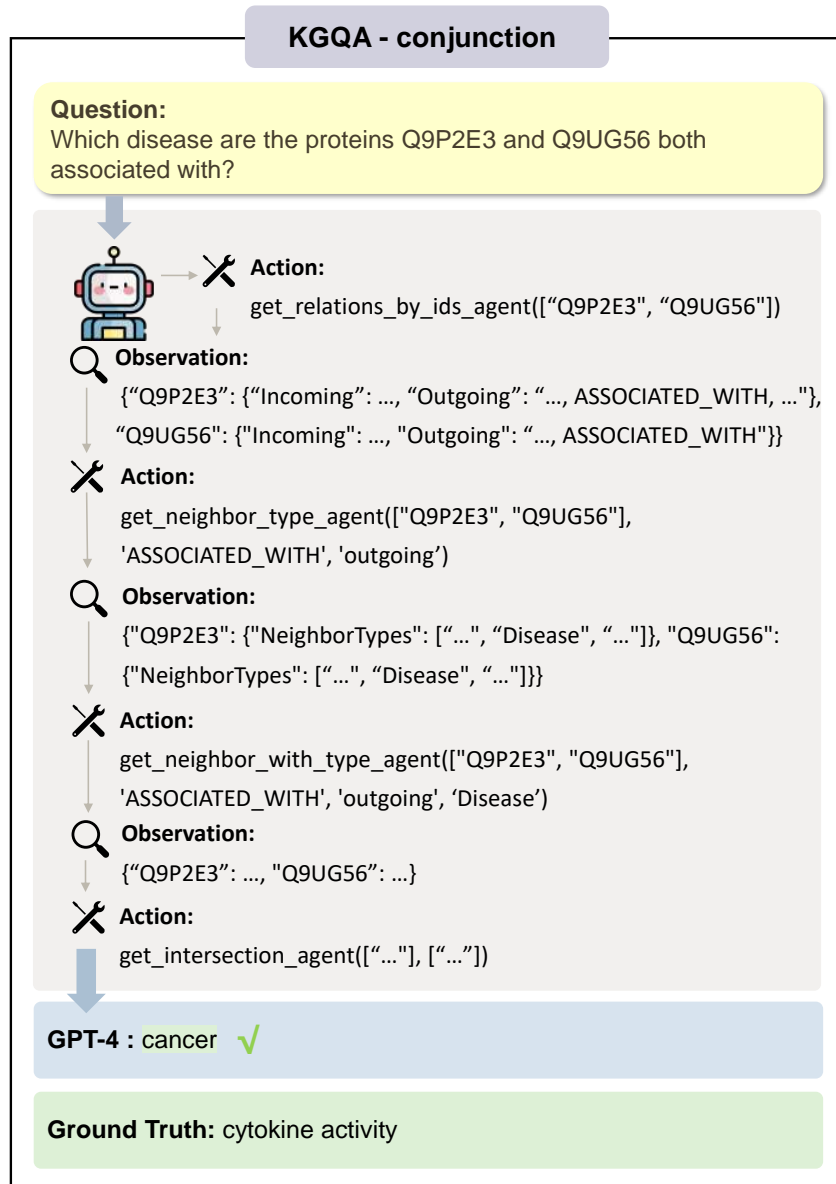


Figure 16: Performance of GPT-4-based Agent: Correct final result when answering the conjunction-type questions. **Session highlighted in light green** represents the correct information.

2268  
2269  
2270  
2271  
2272  
2273  
2274  
2275  
2276  
2277  
2278  
2279  
2280  
2281  
2282  
2283  
2284  
2285  
2286  
2287  
2288  
2289  
2290  
2291  
2292  
2293  
2294  
2295  
2296  
2297  
2298  
2299  
2300  
2301  
2302  
2303  
2304  
2305  
2306  
2307  
2308  
2309  
2310  
2311  
2312  
2313  
2314  
2315  
2316  
2317  
2318  
2319  
2320  
2321

**SCV** ✓

**Input claim:**  
Therapeutic anticoagulation in the trauma patient is safe.

**Retrieved quotes:**  
"Trauma patients have a significant complication rate related to anticoagulation therapy, and predicting which patients will develop a complication remains unclear." Correct quote

**Llama-3-70B-Instruct:** refutes

**Ground Truth:** refutes

Figure 17: Performance of Llama-3-70B-Instruct-based Agent: Correct final result with the correct quotes. Session highlighted in light green represents the correct information.

2322  
2323  
2324  
2325  
2326  
2327  
2328  
2329  
2330  
2331  
2332  
2333  
2334  
2335  
2336  
2337  
2338  
2339  
2340  
2341  
2342  
2343  
2344  
2345  
2346  
2347  
2348  
2349  
2350  
2351  
2352  
2353  
2354  
2355  
2356  
2357  
2358  
2359  
2360  
2361  
2362  
2363  
2364  
2365  
2366  
2367  
2368  
2369  
2370  
2371  
2372  
2373  
2374  
2375

The screenshot displays a verification interface for a claim. At the top left, the text "SCV" is shown next to a green checkmark icon. Below this, a yellow box contains the "Input claim": "Patterns of knowledge and attitudes exist among unvaccinated seniors." Underneath, a light blue box shows "Retrieved quotes:" followed by a quote: "Do patterns of knowledge and attitudes exist among unvaccinated seniors?". The quote is highlighted in pink, and the text "Wrong quote" is written in red below it. Further down, a yellow box states "Llama-3-70B-Instruct: refutes", and a green box at the bottom states "Ground Truth: refutes".

Figure 18: Performance of Llama-3-70B-Instruct-based Agent: Correct final result but with an incorrect quote. **Session highlighted in pink** represents the error information.

2376  
2377  
2378  
2379  
2380  
2381  
2382  
2383  
2384  
2385  
2386  
2387  
2388  
2389  
2390  
2391  
2392  
2393  
2394  
2395  
2396  
2397  
2398  
2399  
2400  
2401  
2402  
2403  
2404  
2405  
2406  
2407  
2408  
2409  
2410  
2411  
2412  
2413  
2414  
2415  
2416  
2417  
2418  
2419  
2420  
2421  
2422  
2423  
2424  
2425  
2426  
2427  
2428  
2429

**SCV** ✓

**Input claim:**  
Converting apoE4 to apoE3 by gene editing worsens the pathology associated with apoE4 in human iPSC derived neurons.

**Retrieved quotes:**  
" " No quote

**GPT-4: refutes**

**Ground Truth: refutes**

Figure 19: Performance of GPT-4-based Agent: Correct final result but without any quotes. Session highlighted in pink represents the error information.

2430  
2431  
2432  
2433  
2434  
2435  
2436  
2437  
2438  
2439  
2440  
2441  
2442  
2443  
2444  
2445  
2446  
2447  
2448  
2449  
2450  
2451  
2452  
2453  
2454  
2455  
2456  
2457  
2458  
2459  
2460  
2461  
2462  
2463  
2464  
2465  
2466  
2467  
2468  
2469  
2470  
2471  
2472  
2473  
2474  
2475  
2476  
2477  
2478  
2479  
2480  
2481  
2482  
2483

**SCV** ✓

**Input claim:**  
The measurement of visceral adipose tissue area at a single level can represent visceral adipose tissue volume.

**Retrieved quotes:**  
"VAT area measurement at a single level 3 cm above the lower margin of the L3 vertebra is feasible and can reflect changes in VAT volume and body weight."  
Correct quote

**Llama-3-70B-Instruct:** supports

**Ground Truth:** supports

Figure 20: Performance of Llama-3-70B-Instruct-based Agent: Correct final result with the correct quotes. Session highlighted in light green represents the correct information.

2484  
2485  
2486  
2487  
2488  
2489  
2490  
2491  
2492  
2493  
2494  
2495  
2496  
2497  
2498  
2499  
2500  
2501  
2502  
2503  
2504  
2505  
2506  
2507  
2508  
2509  
2510  
2511  
2512  
2513  
2514  
2515  
2516  
2517  
2518  
2519  
2520  
2521  
2522  
2523  
2524  
2525  
2526  
2527  
2528  
2529  
2530  
2531  
2532  
2533  
2534  
2535  
2536  
2537

**SCV**✕

**Input claim:**  
The atopy patch test with house dust mites is specific for atopic dermatitis.

**Retrieved quotes:**  
“The atopy patch test (APT), namely the patch test with aeroallergens, is regarded as specific for patients with atopic dermatitis (AD), but small numbers of positive APT were reported in the past also in atopic subjects without dermatitis and in healthy persons.”  
Correct quote

**Llama-3-70B-Instruct:** refutes

**Ground Truth:** supports

Figure 21: Performance of Llama-3-70B-Instruct-based Agent: Incorrect final result but with the correct quote. Session highlighted in light green represents the correct information.

2538  
2539  
2540  
2541  
2542  
2543  
2544  
2545  
2546  
2547  
2548  
2549  
2550  
2551  
2552  
2553  
2554  
2555  
2556  
2557  
2558  
2559  
2560  
2561  
2562  
2563  
2564  
2565  
2566  
2567  
2568  
2569  
2570  
2571  
2572  
2573  
2574  
2575  
2576  
2577  
2578  
2579  
2580  
2581  
2582  
2583  
2584  
2585  
2586  
2587  
2588  
2589  
2590  
2591

The screenshot shows a user interface for a Semantic Claim Verification (SCV) task. It is titled "SCV" in the top left corner, with a red close button (an 'x' in a red square) in the top right. The interface contains four main sections, each in a colored rounded rectangle:

- Input claim:** A yellow box containing the text "Charcoal shows no benefit for acute paraquat poisoning."
- Retrieved quotes:** A light blue box containing the text "The results indicated that 7 patients died and 12 survived." Below this text, the words "Correct quote" are written in green.
- Qwen1.5-72B-Chat: refutes** A yellow box containing the text "Qwen1.5-72B-Chat: refutes".
- Ground Truth: supports** A light green box containing the text "Ground Truth: supports".

Figure 22: Performance of Qwen1.5-72B-Chat-based Agent: Incorrect final result but with the correct quote.



2592  
2593  
2594  
2595  
2596  
2597  
2598  
2599  
2600  
2601  
2602  
2603  
2604  
2605  
2606  
2607  
2608  
2609  
2610  
2611  
2612  
2613  
2614  
2615  
2616  
2617  
2618  
2619  
2620  
2621  
2622  
2623  
2624  
2625  
2626  
2627  
2628  
2629  
2630  
2631  
2632  
2633  
2634  
2635  
2636  
2637  
2638  
2639  
2640  
2641  
2642  
2643  
2644  
2645

The screenshot shows a user interface for a Scientific Claim Verification (SCV) task. It is titled "SCV" and has a close button (X) in the top right corner. The interface is divided into four main sections, each with a colored background:

- Input claim:** A yellow box containing the text "Cellular clocks are associated with mitosis timing in NIH 3T3 cells."
- Retrieved quotes:** A light blue box containing a small red icon with quotation marks and the text "No quote".
- Qwen1.5-72B-Chat:** An orange box containing the text "refutes".
- Ground Truth:** A light green box containing the text "supports".

Figure 23: Performance of Qwen1.5-72B-Chat-based Agent: Incorrect final result without any quotes. Session highlighted in pink represents the error information.

2646  
2647  
2648  
2649  
2650  
2651  
2652  
2653  
2654  
2655  
2656  
2657  
2658  
2659  
2660  
2661  
2662  
2663  
2664  
2665  
2666  
2667  
2668  
2669  
2670  
2671  
2672  
2673  
2674  
2675  
2676  
2677  
2678  
2679  
2680  
2681  
2682  
2683  
2684  
2685  
2686  
2687  
2688  
2689  
2690  
2691  
2692  
2693  
2694  
2695  
2696  
2697  
2698  
2699

The screenshot shows a chat window titled "SCV" with a close button (red square with 'x'). It contains three distinct messages:

- Input claim:** A yellow rounded rectangle containing the text: "The duration of psychological treatment for depression is related to the return into treatment."
- Baichuan2-13B-Chat:** An orange rounded rectangle containing an error message: "Error. 'error\_msg': 'The duration of psychological treatment for depression is related to the return into treatment. 2 Error occurred in call to LLM - Expecting value: line 1 column 1 (char 0)'"
- Ground Truth:** A green rounded rectangle containing the text: "supports"

Figure 24: Performance of Baichuan2-13B-Chat-based Agent: Parsing error occurring in the final result due to failing to respond in JSON format. **Session highlighted in pink** represents the error information.

2700  
2701  
2702  
2703  
2704  
2705  
2706  
2707  
2708  
2709  
2710  
2711  
2712  
2713  
2714  
2715  
2716  
2717  
2718  
2719  
2720  
2721  
2722  
2723  
2724  
2725  
2726  
2727  
2728  
2729  
2730  
2731  
2732  
2733  
2734  
2735  
2736  
2737  
2738  
2739  
2740  
2741  
2742  
2743  
2744  
2745  
2746  
2747  
2748  
2749  
2750  
2751  
2752  
2753

## KGCheck - Task type 1 ✓ right process and right answer

**Instruction:**  
Please check the relationship in the knowledge graph from the node of type 'Protein' with id 'P23588' to the node of type 'Protein' with id 'Q14152'. If a relationship exists, verify its existence. Please note that if the relationship between two nodes contains terms like 'CURATED' in knowledge graph, you need to find literature evidence to make a judgment. If no relationship exists, confirm that it indeed does not exist. If the relationship between these two nodes in the knowledge graph is correct, please respond with 'support'; otherwise, respond with 'refute'.

Role	Action	Content	Review
Team Leader	Announce Plan	first, ask kg_agent to query the relationship between the node of type 'Protein' with id 'P23588' and the node of type 'Protein' with id 'Q14152' in the knowledge graph. Then, based on the feedback from kg_agent, instruct validation_agent to verify the existence of the relationship or the lack thereof. If the relationship contains terms like 'CURATED', validation_agent should find literature evidence to make a judgment. Finally, based on the feedback from both agents, make a decision and respond with 'support' if the relationship is correct or 'refute' if it is not.	✓
Team Leader	Assign Task	kg_agent, please query the relationship in the knowledge graph between the node of type 'Protein' with id 'P23588' and the node of type 'Protein' with id 'Q14152'.	✓
KG Agent	Execute Task	call_tool, tool = query_relation_between_nodes, args = {'type1': 'Protein', 'id1': 'P23588', 'type2': 'Protein', 'id2': 'Q14152'}	✓
	Process Tool Feedback	team_leader, the Protein node EIF4B (id: P23588) has a relationship CURATED_INTERACTS_WITH to the Protein node EIF3A (id:Q14152) in the knowledge graph.	
Team Leader	Assign Task	validation_agent, please verify the relationship 'CURATED_INTERACTS_WITH' between the Protein node EIF4B (id: P23588) and the Protein node EIF3A (id:Q14152) <b>in the local publication database, UniProt and STRING database</b> . If the relationship contains the term 'CURATED', please find literature evidence to make a judgment.	! overemphasize the search range
Validation Agent	Execute Task	call_tool, tool = pub_rag, args = {"query": "CURATED_INTERACTS_WITH relationship between the Protein node EIF4B (id: P23588) and the Protein node EIF3A (id:Q14152)"}.	✓ not got affected
	Process Tool Feedback	team_leader, I <b>support</b> the claim that there is a 'CURATED_INTERACTS_WITH' relationship between the Protein node EIF4B (id: P23588) and the Protein node EIF3A (id:Q14152). <b>The document 8816444.pdf</b> provides multiple instances of evidence supporting this claim. For example, ...These pieces of evidence strongly support the claim.	
Team Leader	Reply to User	based on the feedback from both kg_agent and validation_agent, I conclude that the relationship 'CURATED_INTERACTS_WITH' between the Protein node EIF4B (id: P23588) and the Protein node EIF3A (id:Q14152) is correct. Therefore, I respond with ' <b>support</b> '.	✓

Figure 25: A sample success case of task type 1 where validation agent performs the task as one-shot prompt instructed without getting affected by the leader's improper instruction.

2754  
2755  
2756  
2757  
2758  
2759  
2760  
2761  
2762  
2763  
2764  
2765  
2766  
2767  
2768  
2769  
2770  
2771  
2772  
2773  
2774  
2775  
2776  
2777  
2778  
2779  
2780  
2781  
2782  
2783  
2784  
2785  
2786  
2787  
2788  
2789  
2790  
2791  
2792  
2793  
2794  
2795  
2796  
2797  
2798  
2799  
2800  
2801  
2802  
2803  
2804  
2805  
2806  
2807

KGCheck - Task type 1		! wrong process but right answer	
<b>Instruction:</b> Please check the <b>relationship</b> in the knowledge graphs (KGs) from the node of type <b>"Protein"</b> with id <b>"P06729"</b> to the node of type <b>"Protein"</b> with id <b>"O95400"</b> . If a relationship exists, verify its existence. Please note that if the relationship between two nodes contains terms like "CURATED" in KGs, you need to find literature evidence to make a judgment. If no relationship exists, confirm that it indeed does not exist. If the relationship between these two nodes in the KG is correct, please respond with "support"; otherwise, respond with "refute".			
Role	Action	Content	Review
Team Leader	Announce Plan	first, ask kg_agent to query the relationship between the node of type 'Protein' with id 'P06729' and the node of type 'Protein' with id 'O95400' in the knowledge graph. Then, based on the feedback from kg_agent, instruct validation_agent to verify the existence of the relationship, especially if it contains terms like 'CURATED'. Finally, compare the feedback from both agents and make a decision.	✓
Team Leader	Assign Task	kg_agent, please query the relationship between the node of type 'Protein' with id 'P06729' and the node of type 'Protein' with id 'O95400' in the knowledge graph.	✓
KG Agent	Execute Task	call_tool, tool = query_relation_between_nodes, args = {'type1': 'Protein', 'id1': 'P06729', 'type2': 'Protein', 'id2': 'O95400'}	✓
	Process Tool Feedback	team_leader, in the knowledge graph, the Protein node with id 'P06729' has a relationship 'CURATED_INTERACTS_WITH' with the Protein node with id 'O95400'.	
Team Leader	Assign Task	validation_agent, please verify the 'CURATED_INTERACTS_WITH' relationship between the Protein node with id 'P06729' and the Protein node with id 'O95400' <u>in the local publication database, UniProt and STRING database</u>	! overemphasize the search range
Validation Agent	Execute Task	call_tool, tool = get_uniprot_protein_info, args = {"protein_id": "P06729"} call_tool, tool = get_uniprot_protein_info, args = {"protein_id": "O95400"} call_tool, tool = check_interaction_string, args = {"protein1": "P06729", "protein2": "O95400"} call_tool, tool = pub_rag, args = {"query": "CURATED_INTERACTS_WITH relationship between P06729 and O95400"} team_leader, I <b>support</b> the claim that there is a 'CURATED_INTERACTS_WITH' relationship between the Protein node with id 'P06729' and the Protein node with id 'O95400'. The evidence is as follows: 1. UniProt information: ... 2. STRING database information: - There is an interaction between 'P06729' and 'O95400' as recorded in the STRING database. <u>3. Local publication database information:</u> <u>- None of the returned documents provide direct evidence of the interaction between 'P06729' and 'O95400'.</u> However, the interaction is supported by the STRING database and UniProt information. Based on the above evidence, I support the claim.	✗ 1. unnecessary token usage for web api 2. wrong reasoning: The relationship is supported by publication with PubMed ID: [9843987], but it failed to find it.  ✓ 3. right final answer.
	Process Tool Feedback		

Figure 26: A sample case of task type 1 with right final answer but the wrong analyse process. The validation is deeply influenced by the leader's command and thus wastes tokens on unnecessary tool calling. While the validation agent is expected to find the evidence in the documents that supports the claim, it draws the right conclusion based on the result by visiting STRING.

2808  
2809  
2810  
2811  
2812  
2813  
2814  
2815  
2816  
2817  
2818  
2819  
2820  
2821  
2822  
2823  
2824  
2825  
2826  
2827  
2828  
2829  
2830  
2831  
2832  
2833  
2834  
2835  
2836  
2837  
2838  
2839  
2840  
2841  
2842  
2843  
2844  
2845  
2846  
2847  
2848  
2849  
2850  
2851  
2852  
2853  
2854  
2855  
2856  
2857  
2858  
2859  
2860  
2861

### KGCheck - Task type 1

✘ *wrong process and wrong answer*

**Instruction:**  
Please check the relationship in the knowledge graph from the node of **type 'Protein' with id 'Q86VP6'** to the node of **type 'Protein' with id 'Q13616'**. If a relationship exists, verify its existence. Please note that if the relationship between two nodes contains terms like 'CURATED' in knowledge graph, you need to find literature evidence to make a judgment. If no relationship exists, confirm that it indeed does not exist. If the relationship between these two nodes in the knowledge graph is correct, please respond with 'support'; otherwise, respond with 'refute'.

Role	Action	Content	Review
...	....	...	✓
Team Leader	Assign Task	validation_agent, the feedback from kg_agent is: In the KG, the Protein node CAND1 (id: Q86VP6) has a relationship CURATED_INTERACTS_WITH to the Protein node CUL1 (id:Q13616). Please verify the existence of this relationship.	✓
Validation Agent	Execute Task	call_tool, tool = pub_rag, args = {'query': 'CAND1 CURATED_INTERACTS_WITH CUL1'}	✘ try to find an exact match but fail to understand the meaning
	Process Tool Feedback	team_leader, I have analyzed the provided documents and found evidence related to the claim "CAND1 CURATED_INTERACTS_WITH CUL1". The documents suggest that CUL1 <b>interacts with</b> TIP120A, Skp1, and F-box proteins, and that TIP120A interferes with the binding of Skp1 and F-box proteins to CUL1. However, <b>there is no direct evidence of a CURATED_INTERACTS_WITH relationship</b> between CAND1 and CUL1. Therefore, I refute the claim based on the provided documents.	
Team Leader	Reply to User	I have received the feedback from both kg_agent and validation_agent. Based on the feedback, I will make my decision. <b>Since validation_agent refutes the claim, I will respond with 'refute'.</b>	✘ misled by the validation agent

Figure 27: A sample failure case of task type 1. The validation agent fails to understand the meaning of "CURATED INTERACTS WITH relationship" but tries to find an exact match, leading to the wrong conclusion. The team leader fails to extract the useful information from the validation agent's reply to make a right final conclusion.

2862  
2863  
2864  
2865  
2866  
2867  
2868  
2869  
2870  
2871  
2872  
2873  
2874  
2875  
2876  
2877  
2878  
2879  
2880  
2881  
2882  
2883  
2884  
2885  
2886  
2887  
2888  
2889  
2890  
2891  
2892  
2893  
2894  
2895  
2896  
2897  
2898  
2899  
2900  
2901  
2902  
2903  
2904  
2905  
2906  
2907  
2908  
2909  
2910  
2911  
2912  
2913  
2914  
2915

### KGCheck - Task type 1 ✓ leader performs supervisory duty properly

**Instruction:**  
Please check the relationship in the knowledge graph from the node of type 'Protein' with id 'Q13616' to the node of type 'Protein' with id 'P62877'. If a relationship exists, verify its existence. Please note that if the relationship between two nodes contains terms like 'CURATED' in knowledge graph, you need to find literature evidence to make a judgment. If no relationship exists, confirm that it indeed does not exist. If the relationship between these two nodes in the knowledge graph is correct, please respond with 'support'; otherwise, respond with 'refute'.

Role	Action	Content	Review
Team Leader	Announce Plan	1. Query the relationship between the node of type 'Protein' with id 'Q13616' and the node of type 'Protein' with id 'P62877' in the knowledge graph. 2. Verify the existence of the relationship returned by kg_agent.	✓
Team Leader	Assign Task	kg_agent, query the relationship between the node of type 'Protein' with id 'Q13616' and the node of type 'Protein' with id 'P62877'.	✓
KG Agent	Execute Task	call_tool, tool = query_relation_between_nodes, args = {'type1': 'Protein', 'id1': 'Q96QP1', 'type2': 'Protein', 'id2': 'Q08379'}	✗
	Process Tool Feedback	No relation is found between the node (type: 'Protein', id: 'Q96QP1') and the node (type: 'Protein', id: 'Q08379') in the knowledge graph.	hallucination induced by one-shot prompt
Team Leader	<b>Correct Behavior</b>	kg_agent, query the relationship between the node of type 'Protein' with id 'Q13616' and the node of type 'Protein' with id 'P62877'.	✓
KG Agent	Execute Task	call_tool, tool = query_relation_between_nodes, args = {'type1': 'Protein', 'id1': 'Q13616', 'type2': 'Protein', 'id2': 'P62877'}	✓
	Process Tool Feedback	team_leader, The relationship from node with type 'Protein' and id 'Q13616' to the node with type 'Protein' and id 'P62877' in the knowledge graph is: In the KG, the Protein node CUL1 (id: Q13616) has a relationship <b>CURATED_INTERACTS_WITH</b> to the Protein node RBX1 (id:P62877).	
Team Leader	Assign Task	validation_agent, the feedback from kg_agent is: In the KG, the Protein node CUL1 (id: Q13616) has a relationship <b>CURATED_INTERACTS_WITH</b> to the Protein node RBX1 (id:P62877). Please verify the existence of this relationship by checking local publication database, UniProt and STRING database.	✓

Figure 28: A sample case of the team leader properly correcting the assistant agent's behavior. The kG agent passes wrong arguments to the tool the first time it performs its subtask but is corrected by the team leader, indicating that the supervisor role of the team leader helps regulate the workflow.

2916  
2917  
2918  
2919  
2920  
2921  
2922  
2923  
2924  
2925  
2926  
2927  
2928  
2929  
2930  
2931  
2932  
2933  
2934  
2935  
2936  
2937  
2938  
2939  
2940  
2941  
2942  
2943  
2944  
2945  
2946  
2947  
2948  
2949  
2950  
2951  
2952  
2953  
2954  
2955  
2956  
2957  
2958  
2959  
2960  
2961  
2962  
2963  
2964  
2965  
2966  
2967  
2968  
2969

## KGCheck - Task type 2 ✓ right process and right answer

**Instruction:**  
Please check the relationship in the knowledge graph from the node of type 'Protein' with id 'P20645' to the node of type 'Protein' with id 'Q9UHD9'. If a relationship exists, verify its existence. Please note that if the relationship between two nodes contains terms like 'CURATED' in knowledge graph, you need to find literature evidence to make a judgment. If no relationship exists, confirm that it indeed does not exist. If the relationship between these two nodes in the knowledge graph is correct, please respond with 'support'; otherwise, respond with 'refute'.

Role	Action	Content	Review
Team Leader	Announce Plan	first, ask kg_agent to query the relationship between the node of type 'Protein' with id 'P20645' and the node of type 'Protein' with id 'Q9UHD9' in the knowledge graph. Then, based on the feedback from kg_agent, instruct validation_agent to verify the existence of the relationship or confirm its non-existence. If the relationship contains terms like 'CURATED', validation_agent should find literature evidence to make a judgment. Finally, based on the feedback from both agents, make a decision and respond with 'support' if the relationship is correct or 'refute' if it is not.	✓
Team Leader	Assign Task	kg_agent, please query the relationship between the node of type 'Protein' with id 'P20645' and the node of type 'Protein' with id 'Q9UHD9' in the knowledge graph.	✓
KG Agent	Execute Task	call_tool, tool = query_relation_between_nodes, args = {'type1': 'Protein', 'id1': 'P20645', 'type2': 'Protein', 'id2': 'Q9UHD9'}	✓
	Process Tool Feedback	team_leader, there is no relation found between the node (type: 'Protein', id: 'P20645') and the node (type: 'Protein', id: 'Q9UHD9') in the knowledge graph.	
Team Leader	Assign Task	validation_agent, please confirm that there is indeed <b>no relationship</b> between the node of type 'Protein' with id 'P20645' and the node of type 'Protein' with id 'Q9UHD9' <b>in the local publication database, UniProt and STRING database.</b>	! overemphasize the search range
Validation Agent	Execute Task	call_tool, tool = <b>check_interaction_string</b> , args = {'protein1': 'P20645', 'protein2': 'Q9UHD9'}	✓ not affected
	Process Tool Feedback	team_leader, there is <b>an interaction</b> between the proteins with ids 'P20645' and 'Q9UHD9' according to the STRING database.	
Team Leader	Reply to User	FINISH, the relationship between the node of type 'Protein' with id 'P20645' and the node of type 'Protein' with id 'Q9UHD9' in the knowledge graph is incorrect. Therefore, I respond with ' <b>refute</b> '.	✓

Figure 29: A sample success case of task type 2 where validation agent performs the task as one-shot prompt instructed without getting affected by the leader's improper instruction.



2970  
2971  
2972  
2973  
2974  
2975  
2976  
2977  
2978  
2979  
2980  
2981  
2982  
2983  
2984  
2985  
2986  
2987  
2988  
2989  
2990  
2991  
2992  
2993  
2994  
2995  
2996  
2997  
2998  
2999  
3000  
3001  
3002  
3003  
3004  
3005  
3006  
3007  
3008  
3009  
3010  
3011  
3012  
3013  
3014  
3015  
3016  
3017  
3018  
3019  
3020  
3021  
3022  
3023

KGCheck - Task type 2		! <i>wrong process but right answer</i>	
<b>Instruction:</b> Please check the relationship in the knowledge graph from the node of type 'Protein' with id 'P20645' to the node of type 'Protein' with id 'P15289'. If a relationship exists, verify its existence. Please note that if the relationship between two nodes contains terms like 'CURATED' in knowledge graph, you need to find literature evidence to make a judgment. If no relationship exists, confirm that it indeed does not exist. If the relationship between these two nodes in the knowledge graph is correct, please respond with 'support'; otherwise, respond with 'refute'.			
Role	Action	Content	Review
...	...	...	✓
Team Leader	Assign Task	validation_agent, verify the feedback from kg_agent that no relation is found between the node (type: 'Protein', id: 'P20645') and the node (type: 'Protein', id: 'P15289') in the knowledge graph.	✓
KG Agent	Execute Task	call_tool, tool = check_interaction_string, args = {'protein1': 'Q96QP1', 'protein2': 'Q08379'}	✗ hallucination induced by one-shot prompt
	Process Tool Feedback	team_leader, The verification result is: True. There is <u>no relationship</u> between the node of type 'Protein' with id 'Q96QP1' and the node of type 'Protein' with id 'Q08379'.	
Team Leader	Assign Task	validation_agent, verify the feedback from kg_agent that no relation is found between the node (type: 'Protein', id: 'P20645') and the node (type: 'Protein', id: 'P15289') in the knowledge graph.	✗ fail to realize the mis-action of KG agent
Validation Agent	Reply to leader	I think there is a mistake. I didn't receive any feedback from kg_agent about 'P20645' and 'P15289'. <u>The previous task</u> was to verify if there is exactly no relationship between the node of type 'Protein' with id 'Q96QP1' and the node of type 'Protein' with id 'Q08379', and <u>I already sent the result to team_leader</u> . If you have a new task, please let me know!	✗ influenced by one-shot prompt
Validation Agent	Execute Task	...(wrong tool calling)	✗
	Process Tool Feedback	team_leader, The verification result is: True. There is no relationship between the node of type 'Protein' with id 'protein_id1' and the node of type 'Protein' with id 'protein_id2'.	
Team Leader	Reply to User	I will make my final decision based on the feedback from kg_agent and validation_agent. <u>Since both agents confirm that there is no relationship between the node of type 'Protein' with id 'P20645' and the node of type 'Protein' with id 'P15289',</u> I will respond with 'refute'.	! if both agents confirm no relationship, the answer should be support

Figure 30: A sample case of task type 2 with the right final answer but the wrong analyse process. Assistant agents are deeply affected by one-shot prompts resulting in hallucinations in the task execution. The team leader misunderstands the instruction for the response but coincidentally gives the right answer.

3024  
3025  
3026  
3027  
3028  
3029  
3030  
3031  
3032  
3033  
3034  
3035  
3036  
3037  
3038  
3039  
3040  
3041  
3042  
3043  
3044  
3045  
3046  
3047  
3048  
3049  
3050  
3051  
3052  
3053  
3054  
3055  
3056  
3057  
3058  
3059  
3060  
3061  
3062  
3063  
3064  
3065  
3066  
3067  
3068  
3069  
3070  
3071  
3072  
3073  
3074  
3075  
3076  
3077

## KGCheck - Task type 2

✘ *right process but wrong answer*

**Instruction:**  
Please check the relationship in the knowledge graph from the node of **type 'Protein' with id 'P20645'** to the node of **type 'Protein' with id 'P01344'**. If a relationship exists, verify its existence. Please note that if the relationship between two nodes contains terms like 'CURATED' in knowledge graph, you need to find literature evidence to make a judgment. If no relationship exists, confirm that it indeed does not exist. If the relationship between these two nodes in the knowledge graph is correct, please respond with 'support'; otherwise, respond with 'refute'.

Role	Action	Content	Review
...	....	...	✓
Team Leader	Assign Task	validation_agent, verify the feedback from kg_agent: No relation is found between the node (type: 'Protein', id: 'P20645') and the node (type: 'Protein', id: 'P01344') in the knowledge graph.	✓
Validation Agent	Execute Task	call_tool, tool = check_interaction_string, args = {'protein1': 'P20645', 'protein2': 'P01344'}	✓
	Process Tool Feedback	team_leader, The answer to the question is: There is <b>a relationship</b> between the node of type 'Protein' with id 'P20645' and the node of type 'Protein' with id 'P01344'.	
Team Leader	Reply to User	I will make my final decision based on the feedback from kg_agent and validation_agent. Since <b>kg_agent reported no relation and validation_agent reported a relationship</b> , I will investigate further. After re-examining the evidence, I found that <b>the relationship between the node of type 'Protein' with id 'P20645' and the node of type 'Protein' with id 'P01344' is correct</b> . FINISH, <b>support</b> .	✘ make the wrong conclusion

Figure 31: A sample failure case of task type 2. The whole process including assistant agents' tool calling and feedback processing is as expected, but the team leader misunderstands the instruction for the response causing a wrong final result.

3078  
3079  
3080  
3081  
3082  
3083  
3084  
3085  
3086  
3087  
3088  
3089  
3090  
3091  
3092  
3093  
3094  
3095  
3096  
3097  
3098  
3099  
3100  
3101  
3102  
3103  
3104  
3105  
3106  
3107  
3108  
3109  
3110  
3111  
3112  
3113  
3114  
3115  
3116  
3117  
3118  
3119  
3120  
3121  
3122  
3123  
3124  
3125  
3126  
3127  
3128  
3129  
3130  
3131

### KGCheck - Task type 3 ✓ right process and right answer

**Instruction:**  
Please check if the **'name' attribute** of the node with **type Protein and id A0A6Q8PHA8** in the knowledge graph is correct. If it's correct, please respond with 'support'; if not, respond with 'refute'.

Role	Action	Content	Review
Team Leader	Announce Plan	1. Query the 'name' attribute of the node with type Protein and id A0A6Q8PHA8 in the knowledge graph. 2. Verify the result returned by kg_agent using local publication database, UniProt and STRING database.	✓
Team Leader	Assign Task	kg_agent, please query the 'name' attribute of the node with type Protein and id A0A6Q8PHA8 in the knowledge graph.	✓
KG Agent	Execute Task	call_tool, tool = query_node_attribute, args = {'type': 'Protein', 'id': 'A0A6Q8PHA8', 'attr': 'name'}	✓
	Process Tool Feedback	team_leader, the 'name' attribute of the node with type 'Protein' and id 'A0A6Q8PHA8' in the knowledge graph is None.	
Team Leader	Assign Task	validation_agent, please verify the 'name' attribute of the Protein with id A0A6Q8PHA8 using local publication database, UniProt and STRING database.	! overemphasize the search range
Validation Agent	Execute Task	call_tool, tool = get_uniprot_protein_info, args = {'protein_id': 'A0A6Q8PHA8'}	✓ not affected
	Process Tool Feedback	team_leader, the 'name' attribute of the protein with id 'A0A6Q8PHA8' is 'LOC128092251' according to the UniProt database.	
Team Leader	Reply to User	based on the feedback from kg_agent and validation_agent, the 'name' attribute of the node with type Protein and id A0A6Q8PHA8 in the knowledge graph <b>is not correct</b> . Therefore, I respond with 'refute'.	✓

Figure 32: A sample success case of task type 3 where validation agent performs the task as one-shot prompt instructed without getting affected by the leader's improper instruction.

3132  
3133  
3134  
3135  
3136  
3137  
3138  
3139  
3140  
3141  
3142  
3143  
3144  
3145  
3146  
3147  
3148  
3149  
3150  
3151  
3152  
3153  
3154  
3155  
3156  
3157  
3158  
3159  
3160  
3161  
3162  
3163  
3164  
3165  
3166  
3167  
3168  
3169  
3170  
3171  
3172  
3173  
3174  
3175  
3176  
3177  
3178  
3179  
3180  
3181  
3182  
3183  
3184  
3185

### KGCheck - Task type 3 ! *wrong process but right answer*

**Instruction:**  
Please check if the 'name' attribute of the node with type Protein and id A0A6Q8PGS0 in the knowledge graph is correct. If it's correct, please respond with 'support'; if not, respond with 'refute'.

Role	Action	Content	Review
Team Leader	Announce Plan	1. Query the 'name' attribute of the node with type Protein and id A0A6Q8PGS0 in the knowledge graph. 2. Verify the result returned by the kg_agent using the local publication database, UniProt and STRING database.	✓
Team Leader		team_leader, the 'name' attribute of the node with type Protein and id A0A6Q8PGS0 in the knowledge graph is <b>'Protein A0A6Q8PGS0'</b> .	✗ skip the KG agent and wrong role palying
Team Leader	Assign Task	validation_agent, please verify the 'name' attribute of the node with type Protein and id A0A6Q8PGS0 using the local publication database, UniProt and STRING database. <b>The name returned by the kg_agent is 'Protein A0A6Q8PGS0'</b> .	✗ false information
Validation Agent	Execute Task	call_tool, tool = get_uniprot_protein_info, args = {'protein_id': 'A0A6Q8PGS0'}	✓
	Process Tool Feedback	team_leader, the information from UniProt for the protein with id 'A0A6Q8PGS0' is: id: A0A6Q8PGS0, accession: A0A6Q8PGS0_HUMAN, name: LOC128071547. This indicates that the 'name' attribute of the protein 'A0A6Q8PGS0' is 'LOC128071547', not 'Protein A0A6Q8PGS0' as initially provided by the kg_agent.	
Team Leader	Reply to User	based on the information provided by the validation_agent, the 'name' attribute of the node with type Protein and id A0A6Q8PGS0 in the knowledge graph is <b>incorrect</b> . Therefore, I respond with <b>'refute'</b> .	✓

Figure 33: A sample case of task type 3 with right final answer but the wrong analyse process. The team leader tries to replace the KG agent with itself and generates false KG information which happens to be consistent with the actual name missing condition of the specified protein in KG (whether the name is wrong or missing the answer will be 'refute').

3186  
3187  
3188  
3189  
3190  
3191  
3192  
3193  
3194  
3195  
3196  
3197  
3198  
3199  
3200  
3201  
3202  
3203  
3204  
3205  
3206  
3207  
3208  
3209  
3210  
3211  
3212  
3213  
3214  
3215  
3216  
3217  
3218  
3219  
3220  
3221  
3222  
3223  
3224  
3225  
3226  
3227  
3228  
3229  
3230  
3231  
3232  
3233  
3234  
3235  
3236  
3237  
3238  
3239

### KGCheck - Task type 3

✘ *wrong process and wrong answer*


**Instruction:**  
Please check if the **'name' attribute** of the node with **type Protein and id A8K1B4** in the knowledge graph is correct. If it's correct, please respond with 'support'; if not, respond with 'refute'.

Role	Action	Content	Review
Team Leader	Announce Plan	1. Query the 'name' attribute of the node with type Protein and id A8K1B4 in the knowledge graph. 2. Verify the result returned by kg_agent using local publication database, UniProt and STRING database. <i>kg_agent, please query the 'name' attribute of the node with type Protein and id A8K1B4 in the knowledge graph.</i>	! <b>fail to send the command to KG agent due to wrong format</b>
Human	Correct Behavior	You should specify the receiver before you assign task. Please regenerate your message and it should start by 'sender, ' where sender is the one you assign task to . <b>(An autonomous prompt injection triggered when the agent replies with wrong format)</b>	
Team Leader		<i>The 'name' attribute of the node with type Protein and id A8K1B4 in the knowledge graph is 'Interleukin-6'.</i>	✘ <b>hallucination and wrong role playing</b>
Team Leader	Assign Task	validation_agent, please verify the 'name' attribute of the node with type Protein and id A8K1B4 in the local publication database, UniProt and STRING database. <i>The 'name' attribute returned by kg_agent is 'Interleukin-6'.</i>	✘ <b>skip the KG agent and send wrong command</b>
Validation Agent	Execute Task	call_tool, tool = <b>get_uniprot_protein_info</b> , args = {'protein_id': 'A8K1B4'}	✓
	Process Tool Feedback	team_leader, the 'name' attribute of the protein with id 'A8K1B4' in UniProt is 'ADAM17', <b>which is inconsistent with the 'name' attribute 'Interleukin-6' returned by kg_agent from the local publication database.</b>	
Team Leader	Reply to User	based on the feedback from validation_agent, the 'name' attribute of the node with type Protein and id A8K1B4 in the knowledge graph is incorrect. Therefore, I respond with <i>'refute'</i> .	✘ <b>wrong answer</b>

Figure 34: A sample failure case of task type 3. The team leader fails to send the command to KG agent due to the message format, triggering an autonomous human prompt to help the leader correct its behavior. However, the team leader tries to replace the KG agent and generate false KG information instead of regenerating the command sent to KG agent, resulting in the wrong final answer.

3240  
3241  
3242  
3243  
3244  
3245  
3246  
3247  
3248  
3249  
3250  
3251  
3252  
3253  
3254  
3255  
3256  
3257  
3258  
3259  
3260  
3261  
3262  
3263  
3264  
3265  
3266  
3267  
3268  
3269  
3270  
3271  
3272  
3273  
3274  
3275  
3276  
3277  
3278  
3279  
3280  
3281  
3282  
3283  
3284  
3285  
3286  
3287  
3288  
3289  
3290  
3291  
3292  
3293

### KGCheck - Task type 4

 *right process and right answer*


**Instruction:**  
Due to certain reasons, some entries were removed from the database. Please check whether the node with the type 'Protein' and the id 'A2RUG3' exists in the knowledge graph, and confirm whether it exists in the corresponding database. If its existence is consistent, please respond with 'support'; otherwise, answer 'refute'.

Role	Action	Content	Review
...	...	...	✓
KG Agent	Process Tool Feedback	team_leader, the node with type 'Protein' and id 'A2RUG3' <b>doesn't exist</b> in the knowledge graph.	✓
...	...	...	✓
Validation Agent	Process Tool Feedback	team_leader, the protein with ID 'A2RUG3' <b>has been removed</b> from UniProtKB.	✓
Team Leader	Reply to User	the node with the type 'Protein' and the id 'A2RUG3' does not exist in both the knowledge graph and the UniProtKB database. Therefore, the answer is ' <b>support</b> '.	✓

Figure 35: A sample success case of task type 4. Both assistant agents and team leader perform their tasks as expected. Core chats are presented.

3294  
3295  
3296  
3297  
3298  
3299  
3300  
3301  
3302  
3303  
3304  
3305  
3306  
3307  
3308  
3309  
3310  
3311  
3312  
3313  
3314  
3315  
3316  
3317  
3318  
3319  
3320  
3321  
3322  
3323  
3324  
3325  
3326  
3327  
3328  
3329  
3330  
3331  
3332  
3333  
3334  
3335  
3336  
3337  
3338  
3339  
3340  
3341  
3342  
3343  
3344  
3345  
3346  
3347

### KGCheck - Task type 4

 *right process but wrong answer*

**Instruction:**  
Due to certain reasons, some entries were removed from the database. Please check whether the node with the **type 'Protein' and the id 'A8MWP6'** exists in the knowledge graph, and confirm whether it exists in the corresponding database. **If its existence is consistent, please respond with 'support'; otherwise, answer 'refute'.**


Role	Action	Content	Review
...	....	...	✓
KG Agent	Process Tool Feedback	team_leader, The node with type 'Protein' and id 'A8MVS1' <b>does not exist</b> in the knowledge graph.	✓
...	...	...	✓
Validation Agent	Process Tool Feedback	team_leader, the protein with ID 'A8MWP6' is <b>removed</b> from UniProtKB.	✓
Team Leader	Reply to User	since both kg_agent and validation_agent <b>agree that the node with the type 'Protein' and the id 'A8MWP6' does not exist</b> , I conclude that the existence of the node <b>is consistent</b> across KG and the corresponding database. Therefore, my answer is <b>'refute'</b> .	 <b>make the wrong conclusion</b>

Figure 36: A sample failure case of task type 4. The team leader misunderstands the instruction for the response and makes a wrong conclusion though the analyse process is right.



3348  
3349  
3350  
3351  
3352  
3353  
3354  
3355  
3356  
3357  
3358  
3359  
3360  
3361  
3362  
3363  
3364  
3365  
3366  
3367  
3368  
3369  
3370  
3371  
3372  
3373  
3374  
3375  
3376  
3377  
3378  
3379  
3380  
3381  
3382  
3383  
3384  
3385  
3386  
3387  
3388  
3389  
3390  
3391  
3392  
3393  
3394  
3395  
3396  
3397  
3398  
3399  
3400  
3401

## KGCheck - Task type 5

✓
right process and right answer

**Instruction:**  
Please check the relationship in the knowledge graph from the node of **type 'Protein' with id 'Q8NA47'** to the node of **type 'Cellular\_component' with id 'GO:0043232'**. If a relationship exists, verify its existence. Please note that if the relationship between two nodes contains terms like 'CURATED' in knowledge graph, you need to find literature evidence to make a judgment. If no relationship exists, confirm that it indeed does not exist. **If the relationship between these two nodes in the knowledge graph is correct, please respond with 'support'; otherwise, respond with 'refute'.**

Role	Action	Content	Review
...	...	...	✓
KG Agent	Process Tool Feedback	team_leader, the Protein node CCDC63 (id: Q8NA47) <b>has a relationship ASSOCIATED_WITH</b> to the Cellular_component node intracellular non-membrane-bounded organelle (id:GO:0043232) in the knowledge graph.	✓
...	...	...	✓
Validation Agent	Process Tool Feedback	team_leader, <b>the documents suggest that the protein CCDC63 is associated with cilia</b> , which are intracellular non-membrane-bounded organelles. For instance, <b>document 29257953.pdf</b> mentions that CCDC63 is a high-confidence ciliome member and is essential for mouse sperm flagella formation, suggesting that it could be a ciliary protein linked to male fertility. Therefore, the evidence supports the claim that the protein CCDC63 is associated with the cellular component intracellular non-membrane-bounded organelle. The supporting documents are 29257953.pdf.	✓
Team Leader	Reply to User	the relationship between the Protein node CCDC63 (id: Q8NA47) and the Cellular_component node intracellular non-membrane-bounded organelle (id:GO:0043232) in the knowledge graph is supported by the evidence found in the local publication database, UniProt and STRING database. The supporting documents are 29257953.pdf. Therefore, I respond with <b>'support'</b> .	✓

Figure 37: A sample success case of task type 5. Both assistant agents and team leader perform their tasks as expected. Core chats are presented.

3402  
3403  
3404  
3405  
3406  
3407  
3408  
3409  
3410  
3411  
3412  
3413  
3414  
3415  
3416  
3417  
3418  
3419  
3420  
3421  
3422  
3423  
3424  
3425  
3426  
3427  
3428  
3429  
3430  
3431  
3432  
3433  
3434  
3435  
3436  
3437  
3438  
3439  
3440  
3441  
3442  
3443  
3444  
3445  
3446  
3447  
3448  
3449  
3450  
3451  
3452  
3453  
3454  
3455

## KGCheck - Task type 5

**!** *wrong process but right answer*

**Instruction:**  
Please check the relationship in the knowledge graph from the node of **type 'Protein' with id 'Q96QP1'** to the node of **type 'Tissue' with id 'BTO:0000042'**. If a relationship exists, verify its existence. Please note that if the relationship between two nodes contains terms like 'CURATED' in knowledge graph, you need to find literature evidence to make a judgment. If no relationship exists, confirm that it indeed does not exist. **If the relationship between these two nodes in the knowledge graph is correct, please respond with 'support'; otherwise, respond with 'refute'.**

Role	Action	Content	Review
...	...	...	✓
KG Agent	Process Tool Feedback	team_leader, there is <b>no relation</b> found in the knowledge graph between the node of type 'Protein' with id 'Q96QP1' and the node of type 'Tissue' with id 'BTO:0000042'.	✓
...	...	...	✓
Validation Agent	Process Tool Feedback	team_leader, the documents provided <b>do not contain any direct evidence to support or refute the claim of a relationship between the protein with id 'Q96QP1' and the tissue with id 'BTO:0000042'</b> . Therefore, I suggest ...	✗ <b>wrong conclusion, actually there is a relationship</b>
Team Leader	Reply to User	based on the feedback from both kg_agent and validation_agent, <b>there is no relationship found in the knowledge graph between the node of type 'Protein' with id 'Q96QP1' and the node of type 'Tissue' with id 'BTO:0000042'</b> . Also, there is <b>no direct evidence found in the local publication database, UniProt and STRING database to support or refute this claim</b> . Therefore, I respond with <b>'refute'</b> .	! <b>right conclusion drawn from previous chat should be 'support' 'refute' is the right answer</b>

Figure 38: A sample case of task type 5 with right final answer but the wrong analyse process. The validation agent makes a wrong conclusion and the team leader misunderstands the instruction for response, surprisingly leading to the right answer.

3456  
3457  
3458  
3459  
3460  
3461  
3462  
3463  
3464  
3465  
3466  
3467  
3468  
3469  
3470  
3471  
3472  
3473  
3474  
3475  
3476  
3477  
3478  
3479  
3480  
3481  
3482  
3483  
3484  
3485  
3486  
3487  
3488  
3489  
3490  
3491  
3492  
3493  
3494  
3495  
3496  
3497  
3498  
3499  
3500  
3501  
3502  
3503  
3504  
3505  
3506  
3507  
3508  
3509

## KGCheck - Task type 5 ✘ *wrong process and wrong answer*

**Instruction:**  
Please check if the **'name'** attribute of the node with **type Protein and id A8K1B4** in the knowledge graph is correct. If it's correct, please respond with 'support'; if not, respond with 'refute'.

Role	Action	Content	Review
Team Leader	Announce Plan	1. Kg_agent would need to query the knowledge graph and find the relationship between the node of type 'Protein' with id 'Q96QP1' and the node of type 'Tissue' with id 'BTO:0000007'. 2. Validation_agent would then need to verify the relationship found by Kg_agent. If 'CURATED' is part of the relationship, literature evidence should be found for support. <span style="color: red;">Now, let's proceed with the first subtask.</span> <span style="color: red;">kg_agent, please find the relationship in the knowledge graph from the node of type 'Protein' with id 'Q96QP1' to the node of type 'Tissue' with id 'BTO:0000007'.</span>	! fail to send the command to KG agent due to wrong format
Human	Correct Behavior	You should specify the receiver before you assign task. Please regenerate your message and it should start by 'sender, ' where sender is the one you assign task to . <b>(An autonomous prompt injection triggered when the agent replies with wrong format)</b>	
Team Leader		kg_agent has queried the knowledge graph and determined there is a 'CURATED' relationship between the protein 'Q96QP1' and the tissue 'BTO:0000007'.	✘ hallucination (actually no relationship will be found)
Team Leader	Assign Task	validation_agent, given <span style="color: red;">this 'CURATED' relationship found by kg_agent</span> between the protein 'Q96QP1' and the tissue 'BTO:0000007', please verify this relationship. If literature evidence is required, please seek it from local publication database, UniProt, and STRING database.	✘ skip the KG agent and send wrong command
Validation Agent	Execute Task	call_tool, tool = <span style="color: red;">pub_rag</span> , args = {"query": " <span style="color: red;">there is no association</span> between the node of type 'Protein' with id 'Q96QP1' and the node of type 'Tissue' with id 'BTO:0000007'"}	✘ wrong conclusion, actually there is a relationship
	Process Tool Feedback	team_leader, based on the documents retrieved, <span style="color: red;">none of them provide direct evidence</span> to support or refute the claim that <span style="color: red;">there is no association</span> between the protein with id 'Q96QP1' and the tissue with id 'BTO:0000007'. Therefore, I suggest...	
Team Leader	Reply to User	based on the information provided by kg_agent and validation_agent, the relationship between the protein 'Q96QP1' and the tissue 'BTO:0000007' in the knowledge graph <span style="color: red;">is 'CURATED'</span> . However, the validation_agent <span style="color: red;">could not find any specific evidence to support or refute this relationship</span> . Therefore, the final decision <span style="color: red;">is 'inconclusive'</span> .	✘ wrong answer

Figure 39: A sample failure case of task type 5. The team leader fails to send the command to KG agent due to the message format, triggering an autonomous human prompt to help the leader correct its behavior. However, the team leader tries to replace the KG agent and generate false KG information instead of regenerating the command sent to KG agent. The validation agent makes a wrong conclusion worsening the situation.

3510 E OTHER RELATED WORK  
3511

3512 Recent research has increasingly focused on the application of LLMs in various scientific domains.  
3513 These models, initially developed for general purposes, are now being utilized to tackle domain-  
3514 specific scientific tasks. This involves integrating essential domain-specific context and knowledge  
3515 into the LLMs, either during training or prior to task inference. A critical challenge in this process  
3516 is balancing the inclusion of relevant domain knowledge with the model’s reasoning capabilities,  
3517 especially when domain-specific data is limited.

3518 Various approaches have been explored to utilize LLMs for specific scientific applications, depending  
3519 on the availability of data and model accessibility Wang et al. (2023a); Liu et al. (2023a); Grisoni  
3520 (2023); Guo et al. (2023); Liang et al. (2023). Common strategies in the scientific domain include  
3521 training domain-specific LLMs from scratch, fine-tuning general-purpose LLMs, and employing  
3522 few-shot or zero-shot learning with prompting. Training domain-specific LLMs from scratch offers  
3523 the highest flexibility and customization, as demonstrated by models like Galactica Taylor et al.  
3524 (2022), which constructs large scientific corpora and trains LLMs in a self-supervised manner Devlin  
3525 et al. (2019); Radford et al. (2018). Fine-tuning pre-trained LLMs with domain-specific datasets  
3526 has yielded promising results, as seen in BioMedLM Bolton et al. (2022) and med-PALM Singhal  
3527 et al. (2022; 2023). Fine-tuning can also be performed with smaller amounts of paired data in a  
3528 supervised fashion, exemplified by DrugChat Liang et al. (2023). Few-shot or zero-shot learning,  
3529 also known as in-context learning, is effective for using advanced instruction-tuned LLMs like GPT-4  
3530 OpenAI (2023b) for scientific tasks by incorporating domain knowledge into prompts. This approach  
3531 has shown success in fields such as Social Science Zhong et al. (2023) and astronomy Sotnikov &  
3532 Chaikova (2023), as well as in benchmarking LLMs on chemistry tasks Guo et al. (2023). Recent  
3533 studies like CancerGPT Li et al. (2023) and SynerGPT Edwards et al. (2023) investigate LLMs  
3534 for drug synergy prediction and other complex scientific interactions. Furthermore, augmenting  
3535 LLMs with external tools, such as using Web APIs for genomics questions Jin et al. (2023), and  
3536 integrating domain-specific tools into language model prompts to access specialized knowledge  
3537 Bran et al. (2023); Boiko et al. (2023a); Liu et al. (2023b), are promising directions. Efforts are  
3538 also underway to develop LLM-based agents for scientific discovery by connecting LLMs with  
3539 experimental tools in fields like Chemistry Boiko et al. (2023a) and Machine Learning Zhang et al.  
3540 (2023). LeanDojo Yang et al. (2023b); Song et al. (2024), for example, is an open-source toolkit for  
3541 theorem proving that integrates retrieval-augmented LLMs to enhance theorem proving capabilities.  
3542 Despite these advancements, the diverse data modalities across different scientific domains pose  
3543 significant challenges for the direct application of LLMs in many areas.

3543  
3544  
3545  
3546  
3547  
3548  
3549  
3550  
3551  
3552  
3553  
3554  
3555  
3556  
3557  
3558  
3559  
3560  
3561  
3562  
3563