

HypoChainer: A Collaborative System Combining LLMs and Knowledge Graphs for Hypothesis-Driven Scientific Discovery

Haoran Jiang , Shaohan Shi , Yunjie Yao , Chang Jiang , Quan Li 

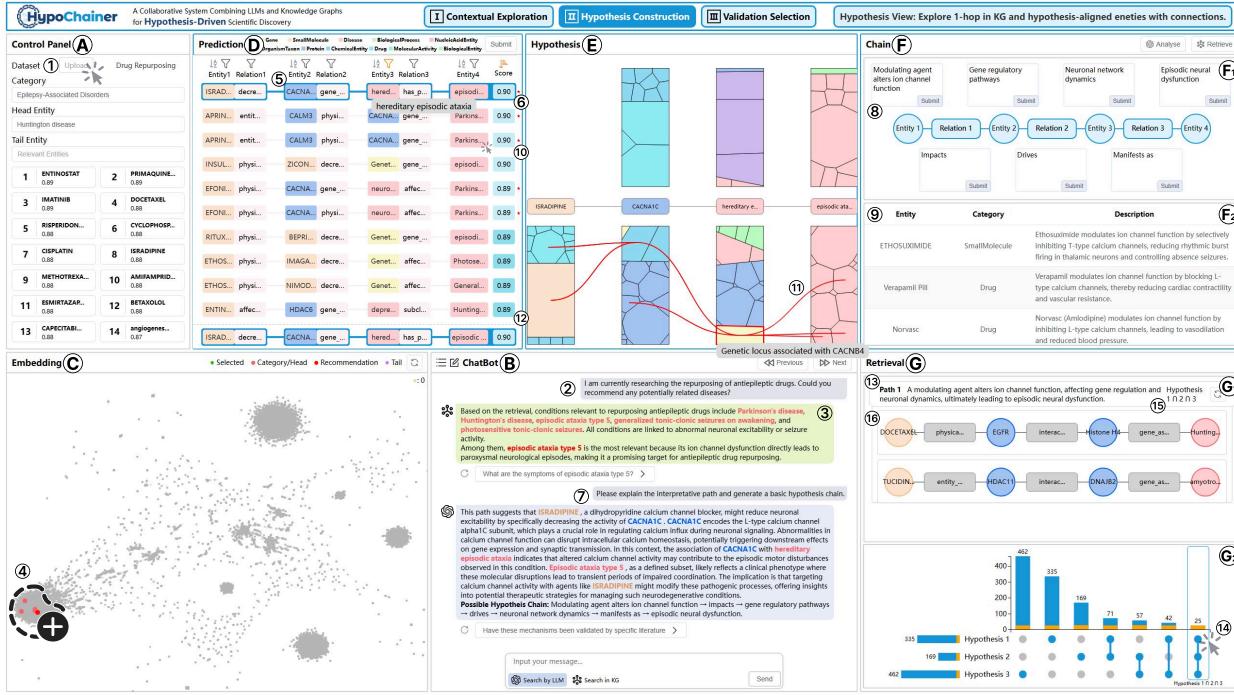


Fig. 1: In the presented case study, the biologist's analytical workflow unfolds as follows: ① Upload Drug Repurposing data. ② Pose a target question to the RAG model. ③ RAG identifies *Episodic Ataxia Type 5* as the top candidate. ④ Lasso the clustered diseases. ⑤ Observe that the *CACNA1C* often appeared among the top-ranked predictions. ⑥ Select the interpretative path of *Episodic Ataxia Type 5*. ⑦ The LLM explains the selected prediction path and generates base hypotheses. ⑧ Construct the hypothesis chain. ⑨ Hypothesis chain is validated for entity alignment. ⑩ KG integration reveals novel paths in the lower layer. ⑪ Observe multiple interpretative paths consistent with the hypothesis. ⑫ Inconsistent predictions for *Huntington's disease* trigger afterwards exploration. ⑬ Expand the retrieval results. ⑭ - ⑮ Filtered predictions confirm relevance to *Parkinson's disease* and *Amyotrophic Lateral Sclerosis (ALS)*. ⑯ Identify overlooked predictions of *Huntington's disease* aligning with the refined hypothesis chain.

Abstract—Modern scientific discovery faces challenges in integrating the rapidly expanding and diverse knowledge required for exploring novel knowledge in biology. While traditional hypothesis-driven research has proven effective, it is constrained by human cognitive limitations, knowledge complexity, and the high costs of trial-and-error experimentation. Deep learning models, particularly graph neural networks (GNNs), have accelerated scientific progress. However, the vast predictions generated make manual selection for experimental validation impractical. Attempts to leverage large language models (LLMs) for filtering predictions and generating novel hypotheses have been impeded by issues such as hallucinations and the lack of structured knowledge grounding, which undermine their reliability. To address these challenges, we propose *HypoChainer*, a collaborative visualization framework that integrates human expertise, LLM-driven reasoning, and knowledge graphs (KGs) to enhance scientific discovery visually. *HypoChainer* operates through three key stages: (1) *Contextual Exploration*: Domain experts employ retrieval-augmented LLMs (RAGs) and visualizations to extract insights and research focuses from vast GNN predictions, supplemented by interactive explanations for in-depth understanding; (2) *Hypothesis Construction*: Experts iteratively explore the KG information relevant to the predictions and hypothesis-aligned entities, gaining knowledge and insights while refining the hypothesis through suggestions from LLMs; and (3) *Validation Selection*: Predictions are prioritized based on the refined hypothesis chains and KG-supported evidence, identifying high-priority candidates for validation. The hypothesis chains are further optimized through visual analytics of the retrieval results. We evaluated the effectiveness of *HypoChainer* in hypothesis construction and scientific discovery through a case study and expert interviews.

Index Terms—Large Language Model, Visual Analytics, Iterative Human-AI Collaboration, Knowledge Graph, Hypothesis Construction

H. Jiang, S. Shi, Y. Jie and Q. Li (corresponding author) are with School of Information Science and Technology, ShanghaiTech University, and Shanghai Engineering Research Center of Intelligent Vision and Imaging, China. E-mail: (fjjianghr2023, shishh2023, yaoyj2024, liquan)@shanghaitech.edu.cn. C. Jiang is with the Shanghai Clinical

Research and Trial Center, Shanghai, China. E-mail: jiangchang@shanghaitech.edu.cn. Haoran Jiang and Shaohan Shi contributed equally to this work.

Manuscript received 1 April 2024; accepted 15 July 2024. Date of Publication 21 October 2024; date of current version 18 July 2024. Digital Object Identifier: xx.xxxx/TVCN.201x.xxxxxxx

1 INTRODUCTION

Modern scientific discovery faces a critical challenge in synthesizing exponentially growing [60], heterogeneous knowledge to drive breakthroughs in data-intensive domains [16], like biomedicine and drug development [77]. Traditional hypothesis-driven research (Fig. 2-[A]) relies on domain experts to manually formulate theories through exhaustive literature reviews, iterative experimentation, and reasoning grounded in specialized knowledge [51]. While this paradigm has yielded significant advances, it's inherently constrained by cognitive limits [7], combinatorial complexity of biological systems [6], and resource-intensive nature of experimental validation (e.g., months-long wet-lab studies for a single hypothesis) [66]. These limitations highlight the need for computational frameworks that enhance human expertise by reducing information overload and reliance on trial-and-error.

Advances in deep learning, particularly graph neural networks (GNNs) and transformer-based models [70, 72], have revolutionized hypothesis prioritization—the process of ranking and selecting the most promising hypotheses—in biological discovery. These methods effectively model biological interactions—such as protein-ligand binding [13], synthetic lethality (SL) relationships¹ [72], and drug candidate prediction [41]—to generate predictive outcomes that guide hypothesis validation. While these tools have greatly streamlined the scientific discovery process, their effectiveness is increasingly constrained by the rapid expansion of biological and biomedical datasets. As model predictions grow in volume [41], manual evaluation and validation by domain experts becomes impractical, creating a critical bottleneck in translating computational insights into scientific breakthroughs. To address this challenge, previous studies have proposed using the comparative analysis of similar predictive outcomes [27] to identify promising yet unverified predictions. Although these methods have shown success in identifying results that align with established mechanisms, they are limited in their ability to detect predictions that lack prior validated analogs. As a result, the identification and validation of novel mechanisms from predictive outcomes remain a significant challenge, hindering the discovery of groundbreaking scientific insights.

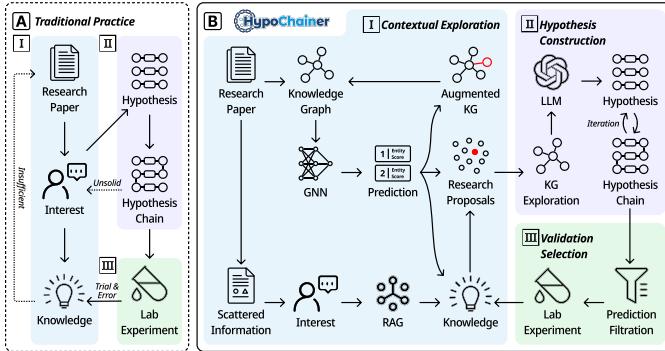


Fig. 2: Comparison of Traditional Practice [A] and HypoChainer Pipeline [B]: Both follow the I Contextual Exploration, II Hypothesis Construction, and III Validation Selection workflow.

Recent efforts [1, 69, 76] have extended to leveraging large language models (LLMs) to analyze and interpret large-scale predictive outcomes. LLMs offer unique advantages in integrating multimodal and heterogeneous data, enabling preliminary reasoning frameworks to address complex biological questions [25, 44]—from hypothesis generation and information retrieval [64] to evidence-based explanations [55]. Yet, as data complexity grows, their propensity for hallucinations or inaccuracies intensifies [50], raising concerns about reliability despite their multi-perspective reasoning capabilities. This tension has spurred interest in grounding LLM outputs in structured knowledge [14, 39] to enhance trustworthiness, where GNNs and knowledge graphs (KGs) have emerged as critical solutions. Particularly, GNNs, renowned for accuracy through structured relational modeling, synergize with

¹Synthetic lethality is a genetic interaction where the simultaneous inhibition of two genes causes specific cell death while inhibiting either gene alone does not.

KGs—which encode precise, standardized relationships—to improve feature representation and prediction robustness [73]. For instance, *KR4SL* [72] integrates KG reasoning with GNNs to predict synthetic lethality, combining semantic relationships and structural patterns to boost both accuracy and explainability. However, challenges persist: textual data alone may inadequately contextualize predictions, while gaps in commonsense knowledge or overly simplistic edge relationships in KGs can limit their comprehensiveness [48].

To leverage the complementary strengths and mitigate the limitations of existing methods (LLMs, KGs, and manual workflows), we focus on a critical challenge: *integrating human expertise, LLM-driven reasoning, and KG-structured knowledge* into a unified hypothesis-driven framework that breaks information silos and streamlines the discovery of novel mechanisms. Through a formative study conducted in collaboration with domain experts, we identified six key design requirements, emphasizing collaborative exploration, interpretability, and iterative hypothesis refinement. Guided by these requirements, we propose a collaborative framework that synergizes human intuition, LLMs, and KGs through the construction of hypothesis chains—structured reasoning paths composed of multiple interrelated hypotheses connected by logical links. The workflow unfolds in three stages: I **Contextual Exploration**: Domain experts raise research questions, prompting a retrieval-augmented LLM (RAG) to surface relevant research objects from GNN predictions. Interactive visualizations and LLM-generated explanations contextualize results, enabling experts to better analyze predictions while addressing gaps in the missing commonsense knowledge and details within structured information. II **Hypothesis Construction**: As experts iteratively analyze predictions, they construct hypothesis chains—semantically linked sequences of insights—supported by LLM-generated refinements and KG relationships. III **Validation Selection**: The workflow filters predictions against refined hypothesis chains, identifying candidates for experimental validation based on alignment with KG-supported evidence. Weak points in the hypothesis chain are further optimized through visual analytics of the retrieval results. To demonstrate the effectiveness of *HypoChainer* (Fig. 2-[B]), we conducted expert interviews and a case study in the field of drug repurposing. In summary, the contributions of this study are as follows:

- We conducted in-depth expert interviews and a thorough literature analysis, identifying six key design requirements for integrating LLMs with structured knowledge in scientific discovery.
- We developed a collaborative framework, *HypoChainer*, linking LLMs and KGs, enabling experts to explore model predictions, construct and refine hypothesis chains, uncover underlying mechanisms, and prioritize validations through interactive visualization.
- We validated *HypoChainer* through a comprehensive case study and expert interviews, demonstrating its effectiveness and generalizability in hypothesis construction and scientific discovery.

2 RELATED WORK

2.1 Hypothesis Generation and Refinement

Hypothesis generation involves proposing new concepts or scientific mechanisms [78], playing a vital role in advancing scientific research [53], which relies heavily on researchers' accumulated knowledge and intuition, introducing limitations and uncertainties. To address these challenges, researchers began leveraging extensive existing knowledge to support hypothesis construction. Early hypothesis generation efforts primarily focused on predicting relationships between concepts, based on the assumption that new ideas emerged from connections with established ones [23, 32]. However, with advancements in language models [9, 74], open-ended idea generation has gained increasing attention [33, 58]. Recent AI-driven hypothesis generation methods employ diverse approaches to conceptualizing research ideas. For instance, *MOOSE* [66] and *IdeaSynth* [52] integrate LLMs into interactive frameworks, facilitating the transition from inspiration to hypothesis construction. Many studies [31, 57, 62] have also combined hypothesis-driven frameworks with visualization-based designs.

Moreover, hypothesis generation is rarely a one-step process, particularly when constructing complex hypothesis chains that require logical

consistency and adherence to fundamental principles. Therefore, hypothesis refinement—driven by feedback and iterative improvements—is equally critical. Methods such as *HypoGeniC* [77] and *MOOSE* [66] emphasize iterative enhancement through feedback mechanisms, including direct responses to hypothesis [4], experimental result evaluations [42, 71], and automated peer-review commentary [40]. Beyond feedback-driven refinements, collaborative hypothesis generation has also gained traction, leading to the development of multi-agent systems [16, 45]. For instance, *VIRSCI* [59] optimizes hypothesis construction by customizing knowledge for each agent, while *Nova* [24] incorporates outputs from other research efforts to refine hypothesis generation. However, such multi-agent frameworks also introduce challenges, including hallucinations [50], inaccuracies, and errors in agent-generated outputs. If left unchecked, these errors can propagate through the reasoning process, leading to misleading or incorrect conclusions. Recently, a multi-agent system based on *Gemini 2.0* was proposed [17], designed to generate and refine novel research hypotheses through a “generate-debate-evolve” framework, utilizing test-time compute to improve hypothesis quality. However, this system does not incorporate constraints from structured knowledge sources and experts still remain limited to guiding AI in knowledge discovery, lacking intuitive visualizations to independently explore and uncover insights.

To address these challenges, we propose combining LLM-driven hypothesis generation with AI-expert collaborative optimization process. By enabling timely human intervention and error correction, it ensures both the logical consistency of the hypothesis and the effective utilization of LLM’s vast knowledge and reasoning capabilities, ultimately enhancing hypothesis construction and iterative refinement.

2.2 Graph-Structured Knowledge Exploration

With the continuous expansion and enrichment of structured knowledge, increasingly complex graphs have emerged across various domains. These graphs encode rich information about entities and their relationships, facilitating reasoning and novel knowledge prediction through propagation. However, as graph data rapidly expands, effectively identifying patterns within large-scale graphs and exploring vast prediction results within the same large-scale KG has become increasingly challenging. Researchers have addressed these issues by developing various visualizations, particularly in biological networks [27, 65], neural networks [29], and social networks [10]. For example, Paley et al. [46] used force-directed layouts and clustering to reveal key biological pathways and drug targets in complex datasets, facilitating discoveries in molecular interactions and metabolic pathways that drive advancements in drug discovery and synthetic biology.

While these visualization methods have demonstrated significant effectiveness, the challenges of graph exploration and information retrieval extend beyond static representations. Many interactive approaches [75] have been developed to tackle these challenges, such as *Biolinker* [13], an interactive visualization system that supports bottom-up exploration of complex protein interaction networks. However, when analyzing predictive model outputs, researchers must navigate large volumes of predictions and interpret extensive graph structures—often with subtle yet critical distinctions. These challenges are further exacerbated when investigating intricate and extended chains of interpretative paths and patterns, making graph analysis increasingly demanding.

Building upon previous research in graph-structured knowledge exploration and leveraging the analytical reasoning capabilities of LLMs, this study integrates graph exploration with text-based prompts. We employ a hierarchical layout to organize one-hop entities along prediction paths and hypothesis-aligned entities—defined as those complying with or related to the given hypothesis descriptions—where hypotheses are initially generated by the LLM and iteratively refined by experts, using Voronoi treemaps [5], enhancing spatial efficiency and clarity. Simultaneously, existing edges within the KG are visualized between entities to highlight structural relationships between predictions and entities from different sources. This dual approach enables users to explore relationships among structurally related and hypothesis-aligned entities, evaluate the coherence between model predictions and proposed hypotheses, and uncover insights through comparative analysis.

2.3 Collaboration Between Human, LLM, and KG

As AI advances and the modes of interaction diversify, the dynamics of both human-AI and AI-to-AI cooperation are continuously evolving. Within the framework of *Collaboration Between Human, LLMs, and KG*, the pairwise interactions [2, 30] can be categorized into three key relationships: **Human-LLM**, **Human-KG**, and **LLM-KG**.

Human-LLM. The interaction between humans and LLMs forms a continuous learning loop. In this collaborative paradigm, humans serve as both users and critical coordinators, guiding context-specific outputs through structured mechanisms. Additionally, reinforcement learning [54] from human feedback further refines LLM outputs to align with ethical considerations and domain-specific objectives [30].

Human-KG. Humans engage with KGs to retrieve, structure, and refine information through queries [49], natural language [56], and visualization tools [3]. By leveraging KGs as structured repositories of factual knowledge, users can enhance information accessibility and explainability in AI applications. For instance, *DrugExplorer* [65] employs explainable AI (XAI) techniques to interpret GNN-based predictions and refine biomedical KGs, while *SLInterpreter* [27] introduces an iterative human-AI collaboration system for SL prediction, which enables experts to enhance KG interpretability and align AI outputs with expertise via metapath strategies and iterative path refinement.

LLM-KG. The bidirectional integration between LLMs and KGs unfolds through two key paradigms: *LLM-enhanced KG* [2] and *KG-enhanced LLM* [47]. LLMs enhance KGs by automating *entity extraction* via prompt-driven mining [28], refining *entity parsing and matching* through synthetic labeled data generation [63], and improving *link prediction* with joint text-KG embeddings [67]. Conversely, KGs enrich LLMs by injecting structured knowledge during *pre-training* [21], optimizing *fine-tuning* with knowledge-aware objectives [37], enhancing *retrieval* through hybrid neural-symbolic architectures [26], and refining *prompt-based reasoning* via KG-guided subgraph extraction and logic-aware chained inference [11].

Recent studies have begun exploring the Collaboration Between Human, LLMs, and KG, such as *KNOWNET* [68], which combines LLMs with KGs to enhance the quality and efficiency of human-AI interaction through structured knowledge validation and iterative query refinement. While integrating LLMs with KGs can enhance knowledge coverage and reasoning, existing approaches often fail to apply LLMs’ reasoning capabilities to KGs directly and underestimate the role of human experts in decision-making, limiting collaborative efficacy. To address this, we propose a hypothesis validation framework that harnesses the complementary strengths of LLMs, KGs, and human experts. By leveraging each component’s strengths and mitigating limitations, our framework promotes transparent, controllable collaboration, enhancing hypothesis validation and discovery outcomes.

3 FORMATIVE STUDY

In this study, we aimed to understand researchers’ challenges, expectations, and requirements for collaborating with AI across various domains comprehensively. To achieve this, semi-structured interviews (see Appendix E for details) with institutional IRB approval were conducted with two researchers specializing in cancer therapy (**E1-E2**), and two researchers focusing on drug research (**E3-E4**) (Mean Age = 39.75, SD = 5.4, 2 males, 2 females). Among them, **E1-E2** specialize in screening SL pairs using the *CRISPR/Cas9* technique, focusing on thyroid cancer (**E1**) and breast cancer (**E2**), respectively. **E3-E4** are engaged in drug research, with **E3** focusing on drug repurposing, while **E4** specializes in discovering new drug targets through pharmacogenomics. Each participant has significant research experience in their respective fields and expertise in AI-assisted prediction. Through inductive coding and thematic analysis [12], we extracted valuable insights into the challenges faced by domain experts and summarized the design requirements. Each interview session lasted approximately 45 minutes.

3.1 Experts’ Conventional Practices

E3, specializing in drug repurposing, outlined their conventional approach to mechanistic investigations [20]. The process begins with

selecting a research domain, followed by an extensive review of background literature to establish foundational knowledge. Combining this information and their domain expertise, they formulate initial hypotheses through logical reasoning, which are then tested experimentally via methods like high-throughput sequencing. However, **E3** highlighted limitations in this traditional workflow: *Research directions are often restricted to familiar fields to reduce uncertainty, and hypothesis development requires laborious literature reviews to balance novelty with plausibility.* This is particularly challenging in less familiar domains, where limited expertise can undermine confidence in hypothesis design. Additionally, over-relying on trial-and-error experiments without theoretical grounding wastes resources and limits progress.

E1 noted that the similar workflows are also applied to SL mechanistic studies, but with one critical distinction: Wet-lab experiments in SL research take around six months to complete, dramatically increasing the cost of trial-and-error. This necessitates rigorous hypothesis evaluation and meticulous experimental planning prior to testing. To facilitate this process, researchers now employ GNN models to predict potential SL pairs. Predictions are filtered using two criteria: (1) the interpretability of interpretative paths² and (2) the model's prediction confidence scores. While this approach reduces reliance on purely familiarity-driven hypothesis, significant challenges remain: *Human intervention is still required to evaluate predictions, and generating truly novel hypotheses continues to demand substantial creativity and domain insight.* Given the large volume of predicted results, the experts have explored the use of LLMs for direct reasoning and filtering [55, 64]. However, when attempting to integrate LLMs into their research, they found it challenging to effectively incorporate local data (e.g., KG). Moreover, their practical use revealed that LLMs often struggle to maintain consistency and coherence in handling complex hypotheses.

3.2 Experts' Concerns and Expectations

Initially, domain experts expressed concerns that structured knowledge—often represented as triplets, the basic units of KGs describing relationships between entities (e.g., *Primaquine-Affects-IKBKG*)—lacks commonsense details and explanatory text. They also noted that the granularity of edge relationships is sometimes too coarse, making it difficult to interpret different paths. As one expert pointed out: “... some GNNs do give us interpretative paths, but honestly, the information these paths contain is pretty limited. You know, even when different groups of entities are linked by the same label, the actual meaning can be very different. [Point at CYLC1 and DAB1] They are both labeled as being involved in cell differentiation, but I'd bet their actual roles are pretty different. All these missing details make interpretative paths hard to understand and, honestly, it just makes me not trust them as much.” This reveals the first challenge: **C1. Insufficient detail and inadequate granularity in interpretative paths.** To address the issue of limited detail and insufficient granularity in interpretative paths [**C1**], it is crucial to provide reliable and comprehensive information. This leads to our first design requirement: **DR1. Enriching interpretative paths with reliable and sufficient detail.**

Experts have also expressed concerns about the potential inadequacy of relevant information for hypothesis construction. When formulating hypotheses, domain experts often require supplementary knowledge to integrate diverse information and generate novel insights. Acquiring such knowledge typically involves extensive literature reviews, which can be burdensome and prone to overlooking critical insights. As **E1** emphasized, “... When aiming to construct novel hypotheses, relying solely on our existing knowledge is often insufficient. We need to gather as much relevant information as possible to inspire new ideas and facilitate reasoning. However, this process is time-consuming and mentally demanding...Sometimes, extensive reading may not always help us grasp key points, and prolonged reading sessions can be distracting.” This highlights the second challenge: **C2. Inefficiency in acquiring relevant knowledge for novel hypothesis construction.** Building on this, experts further emphasized the need for more effective

²The sequence of entities and edges connecting a head entity (e.g., a gene) to a predicted SL partner via biologically meaningful relationships.

methods, such as KG-based information retrieval, semantic search, and recommendation to aid information retrieval, ultimately facilitating hypothesis construction. This leads to the second design requirement: **DR2. Providing efficient methods to obtain relevant information for hypothesis construction.**

Managing large-scale model predictions presents another significant challenge. The most reliable method currently involves domain experts manually analyzing and filtering predictions to identify novel research focuses or by some visualization methods [57]. However, as the volume of predictions increases, this manual or visual-assisted approach becomes unsustainable, placing a heavy burden on experts and potentially affecting their judgment. **E2** illustrated this difficulty: “*Prediction models are pretty impressive—they do get some things right. But with so many results, it's like going from finding a needle in the ocean to a needle in a pond. Trying to sort through everything on our own without a clear direction is still overwhelming and just not realistic.*” This highlights the third challenge: **C3. Impracticality of manually filtering and analyzing large-scale predictions.** In light of the above challenge, **E1** highlighted, “*It's practically infeasible and really limits how much we can focus. Rule-based methods? Yeah, they can simplify things, but they'll definitely miss a lot of results that don't fit perfectly. We really need something smarter and more intuitive to quickly zero in on the results that actually matter.*” This insight leads us to our third design requirement: **DR3. Streamlining and optimizing evaluation and insight extraction in massive predictions.**

Experts also noted that while existing tools can help uncover simple patterns, they often reinforce known results rather than discovering novel insights. As **E2** explained: “*...Some tools can help us with simple patterns, but they just find results that look like what we already know. The problem is, if a model predicts something different—even if it's right—we might not even notice. That keeps us stuck in a loop, always going for the easy, obvious predictions. Suppose we really want to dig deeper and uncover complex mechanisms, then it's not just about finding similar predictions...We need better ways to build logical connections with hypotheses through observations, refine them, and search for patterns that actually lead us to new insights that we might've missed before.*” This reveals the fourth challenge: **C4. Lack of tools to construct hypothesis chains for exploring complex mechanisms.** As **E3** noted, “*When we're building a hypothesis, every step needs to make sense. If one part doesn't, the [whole] thing falls apart. That's why we need something to help us catch [any] flaws. And expanding the hypothesis is tricky too—it can feel pretty random and really depends on what we already know. A tool that suggests reasonable hypotheses and explains them? Yeah, I'd definitely be interested in trying that.*” This leads to the fourth design requirement: **DR4. Supporting interactive construction and continuous optimization of hypothesis chains.**

Even with tools for constructing hypothesis chains, experts found it difficult to retrieve relevant predictions that align with their hypotheses. They emphasized the need for better integration between model predictions and KG to refine hypotheses. **E4** explained, “*If we construct a hypothesis chain based on certain observations and our expertise but are unable to retrieve similar predictions or relevant evidence from the KG, it does not seem to help us come up with more solid hypotheses or further gain new insights*”. This highlights the fifth challenge: **C5. Absence of retrieval methods for hypothesis-aligned predictions.** As **E2** and **E4** noted, “*It's exciting to find known paths or predictions closely aligning with our hypotheses.*” However, designing search logic or rules for each hypothesis proves difficult, they added, “*and on top of that, it is impossible to experimentally validate every emerging hypothesis.*” This highlights the fifth design requirement: **DR5. Providing effective retrieval methods for hypothesis-aligned predictions.**

Experts have also raised concerns about the information flow between LLMs, KGs, and domain experts, particularly based on their prior experiences with LLM integration. Despite efforts to bridge these components, significant challenges and barriers remain in ensuring that information is effectively transmitted and utilized. As **E3** noted after the discussion, “*I'm hoping LLMs and KGs can better share information with each other. Right now, when I use an LLM with a KG, the LLM often relies too much on external knowledge and doesn't fully use the KG*

information unless I explicitly highlight the details. At the same time, the KG struggles to recognize entities when the LLM describes them differently, making the whole interaction awkward and fragmented. I want to make this smoother; improve how they understand each other, and help them combine their strengths to reach stronger conclusions.” This highlights the last challenge: **C6. Information transmission barriers between LLMs, KGs, and domain experts**. Accordingly, we introduce the last design requirement: **DR6. Ensuring seamless and lossless transmission of information between LLMs, KGs, and domain experts**. This requirement aims to consistently maintain the quality of information, ensuring that when accurate information is transmitted between different parties, its completeness and precision are always preserved. Reflecting the perspectives of **E3**, “*Sure, we obviously want the info to be reliable when it’s passed between the LLM, KG, and us. If there’s any mix-up or something gets left out, it’s definitely gonna mess with our judgments. I think this is really the key to making sure we can work together smoothly.*”

4 HYPOTAINER

To ensure a seamless flow of information among different parties and enable experts to efficiently explore the system, the pipeline has been meticulously designed and iteratively refined. It comprises three key modules: **Contextual Exploration**, **Hypothesis Construction**, and **Validation Selection**, aligning with the conventional workflow while improving overall coherence, rationality, and efficiency (Fig. 2) (**DR6**). This section first provides a detailed overview of the backend algorithms implemented, followed by an introduction to the system’s visual design. Finally, a walkthrough case study is presented to offer a concrete and comprehensible introduction to the entire pipeline.

4.1 Data and Backend Engine

We provide an overview of the data utilized in our approach. To assess the effectiveness and generalizability of our method, we employ two types of biological data: *Drug Repurposing* and *Cancer Research*.

Drug Repurposing. A large-scale biomedical knowledge graph [41] was utilized, derived from *RTX-KG2c* (*v2.7.3*), which integrates data from 70 public sources, with 6.4 million entities and 39.3 million edges. The graph is standardized using the Biolink model [61]. To tailor it for drug repurposing, the dataset was refined into a streamlined graph with 3,659,165 entities and 18,291,237 edges. Key data sources include *MyChem*, *SemMedDB*, *NDF-RT*, and *RepoDB*, providing both positive (indications) and negative (contraindications/no-effect) samples. Furthermore, 472 drug-disease pairs are leveraged from *DrugMechDB* for external validation, enhancing the dataset’s reliability and applicability.

KGML-xDTD [41] is a drug repurposing prediction and mechanism of action (MOA) inference model, integrating graph-based learning with reinforcement learning. The module formulates the task as a link prediction problem, which employs *GraphSAGE*, to generate node embeddings by leveraging structural and attribute information from *PubMedBERT*. The MOA module identifies biologically plausible MOA paths using an adversarial actor-critic reinforcement model, which is trained using curated demonstration paths. *KGML-xDTD* enhances prediction (Appendix Tab. 1) while maintaining interpretability.

Cancer Research. This dataset [72] focuses on SL relationships, where the simultaneous inactivation of a gene pair leads to the death of a specific cancer cell. The cancer research knowledge graph is constructed using data from two primary repositories: *SynLethDB*, which catalogs validated gene SL interactions, and *ProteinKG25*, a biomedical knowledge dataset containing information on gene functions, pathways, and biological processes. The resulting KG comprises 42,547 entities, 33 edge types, and a total of 396,619 triplets.

For SL gene pair prediction, we utilize the *KR4SL* model [72], which follows an encoder-decoder architecture and employs a GNN-based approach with a heterogeneous KG. The model predicts SL pairs by tracing relational paths, assigning weights to these paths to capture the strength of SL interactions and uncover potential relationships between unconnected genes. Candidates are ranked based on their likelihood of forming an SL relationship. Each prediction is accompanied by a

three-hop interpretative path. The model achieves a precision of 59% (Appendix Tab. 2), outperforming comparable models.

4.2 Frontend Interface

Working alongside domain experts, we have designed a frontend interface to support scientific discovery workflows through a set of dedicated modules, namely *Control Panel*, *Embedding View*, *Chatbot*, *Prediction View*, *Hypothesis View*, *Chain View*, and *Retrieval View*.

Control Panel. The *Control Panel* (Fig. 1-(A)) includes three search boxes to streamline prediction selection across hierarchical levels. The *Category* search box displays categories to assist in selecting specific types of *Head Entities* or narrowing the scope of subsequent searches, while the *Head* search box leverages an auto-complete function for direct *Head Entity* queries. Upon selecting a *Head Entity*, the *Tail Entity* table positioned below the search boxes automatically populates with the top 50 predicted *Tail Entities*, sorted in ascending order by rank and annotated with their *Name*, *Score*, and *Rank*. However, given that not all datasets contain predefined categories, users can define new categories encompassing relevant entities through the *ChatBot* during the **Contextual Exploration** phase. Each query is logged in the *Category* search box for convenient reuse in subsequent explorations.

Embedding View. The *Embedding View* (Fig. 1-(C)) displays cluster summaries derived from entity features, with methodologies tailored to distinct datasets. For instance, dimensionality reduction for gene data incorporates gene descriptions, nucleotide sequences, and other pertinent information, whereas drug data reduction utilizes attributes like drug descriptions, indications, and mechanisms of action through UMAP [43] (Comparison in Appendix). This aids domain experts in analyzing prediction patterns, offering an intuitive framework to identify research focuses within large-scale predictions (**DR3**). Additionally, the *Embedding View* interacts with the *ChatBot* during the **Contextual Exploration** phase: when the RAG suggests entities relevant to the query, *Embedding View* highlights them for *lasso* selection.

ChatBot. The *ChatBot* (Fig. 1-(B)) serves as the central hub of the system, orchestrating seamless interactions across all interface components by integrating reasoning models ((LLM) and retrieval mechanisms, which is based on *LightRAG* [19] for its lower cost and faster response compared to *GraphRAG* [15] (RAG), to synchronize KG data, AI-generated insights, and expert input (**DR6**). The interface ensures intuitive usability through features such as a left *history* section that logs and summarizes the dialogue, a bottom *input box* with retrieval mode selector (Search in KG & Search by LLM), and upper-right *<< Previous* & *>> Next* buttons that activate phase-specific LLMs. For instance, during the *Contextual Exploration*, RAG and the entity filtering and recommendation are invoked to ensure that the information comes solely from the local database, thus minimizing the risk of recommending entities not present in the KG or generating hallucinations. The chatbox dynamically adapts its color scheme to signal active retrieval modes, serving as a reminder to domain experts of the current model in use. Furthermore, the responses generated by the LLM include suggestions for further exploration, which are based on neighboring entities in the KG or relevant information related to the topic, and are displayed beneath each chatbox. Also, entities that appear in both the *ChatBot* and the corresponding view are color-matched to ensure consistency, helping experts quickly locate and associate relevant entities.

Prediction View. The *Prediction View* (Fig. 1-(D)) is designed to empower users in distilling valuable insights from a substantial volume of predictions (**DR3**), enabling targeted identification of results aligned with their research objectives. Building on the foundational principle of *LineUp* [18], this view adopts a tabular visualization that balances domain experts’ familiarity with advanced analytical capabilities. The design optimizes spatial efficiency by binding entities to their subsequent edge relationships, compressing horizontal layouts without sacrificing contextual continuity. To streamline exploration, the system provides two key features: auto-completion-enabled filtering, which supports multi-hop queries across entities and edge relations, and dynamic highlighting mechanisms that enhance analytical precision. When users hover over entities, the system highlights the full corresponding interpretative path (Fig. 1-(E)), aiding the exploration of

connection patterns. A persistent anchoring bar at the bottom of the view displays the hovered path and fixes the selected one, reducing visual tracking effort and supporting path comparison. Once a hypothesis chain is submitted for analysis in the *Chain view*, hypothesis-aligned predictions are marked by ★ (Fig. 1-⑩) to enable users to quantify the consistency between theoretical assumptions and empirical outcomes. Additionally, experts are able to perform single-column sorting based on path scores and edge weights. These features collectively enable domain experts to iteratively refine insights—from macro-level pattern discovery to granular path analysis—by systematically prioritizing paths through sortable confidence scores, cross-validating hypotheses against prediction, and exporting candidate paths for validation.

Hypothesis View. The *Hypothesis View* (Fig. 1-⑤) presents KG entities aligned with a selected prediction path (**DR2**), enrich structured information (**DR1**), and assist domain experts in hypothesis refinement through iterative optimization (**DR4**). The view dynamically retrieves entities through RAG and organize directly connected 1-hop neighbors in KG (Appendix Fig. 1-①) and hypothesis-aligned entities (Appendix Fig. 1-②) retrieved from the generated hypotheses and the subsequent iterative refined hypotheses in a hierarchical architecture. To optimize layout clarification and scalability for large biomedical datasets, entities within each layer are arranged using a Voronoi treemap [5], a technique proven effective in reducing visual clutter while preserving spatial efficiency, making it particularly suitable for visualizing large-scale knowledge graphs [27]. Entities of the same type are clustered within layers to highlight structural patterns and simplify analysis. When users hover over a node in the Voronoi treemap, the node and its connected edges are highlighted in red, directing attention to relevant relationships.

The view further dynamically establishes plausible KG-derived links between entities across adjacent layers. This approach enables domain experts to progressively and intuitively analyze potential connections between 1-hop entities on the predicted path and hypothesis-aligned semantic matches, fostering a structured understanding of both KG proximity and semantic relevance. By visualizing potential connection patterns, experts can assess whether these connections align with similar underlying hypotheses, thereby refining validation processes and guiding the construction of logically coherent hypothesis chains.

Chain View. The *Chain View* supports domain experts in systematically constructing, previewing, and refining hypotheses (**DR4**). It facilitates the iterative analysis and optimization of individual hypotheses through integration with LLMs, enabling users to chain hypotheses into a cohesive structure for subsequent retrieval and validation.

The *Input Area* (Fig. 1-⑥) enables hypothesis chain construction and review. Each hypothesis node features a text area above for hypothesis description, with additional fields below to define relationships between entities. Clicking the lower right button triggers targeted RAG retrieval to validate the hypothesis against the KG, displaying relevant entities in the *Entity Preview* (Fig. 1-⑦), including entity names, types, and descriptions of how they align with the hypothesis, ordered by their degree of alignment, which is judged by the RAG through a systematically designed prompt. This preview facilitates rapid hypothesis evaluation and refinement without resource-intensive full-scale RAG retrievals. At the top-right corner, the button leverages LLM capabilities to evaluate the hypothesis chain for logical coherence, identifying potential inconsistencies or optimization opportunities. Once validated, the button initiates the formal retrieval process for downstream tasks. To clarify functionality, the button (linked to KG-grounded RAG retrieval) and the button (LLM-driven reasoning) are differentiated through the icons in *ChatBot*.

Retrieval View. The *Retrieval View* provides an overview of all hypothesis chain retrievals and corresponding results. It highlights the quantities of retrieval outcomes aligned with the hypothesis chain at varying matching levels, alongside detailed retrieved results (**DR5**).

The *Retrieval List* (Fig. 1-⑧) catalogs hypothesis chains and their retrieval outcomes, organized into collapsible records. Below, the *UpSet Plot* (Fig. 1-⑨) [35] visualizes the degree to which retrieval results for the currently selected hypothesis chain satisfy the hypothesis, assisting experts in identifying which hypotheses require refinement or commonly appear in predictions. Specifically, the *UpSet Plot* includes three

rows (Fig. 3-①) representing the three hypotheses within the chain, with an adjacent bar chart (Fig. 3-②) showing the count of triplets satisfying individual hypotheses. A central dot matrix and bar chart (Fig. 3-③) represent combinatorial hypothesis satisfaction. For example, the fifth column highlights triplets fulfilling both Hypothesis 2 and 3 in the hypothesis chain. Hovering over a row or column (Fig. 3-④) dynamically highlights intersecting sets and displays their proportional contributions within the corresponding hypothesis (Fig. 3-⑤). Clicking bars filters the *Retrieval List* accordingly, with intersections indicated at the upper-right corner (Fig. 1-⑯). The button resets the view to the default display of complete retrieval results.

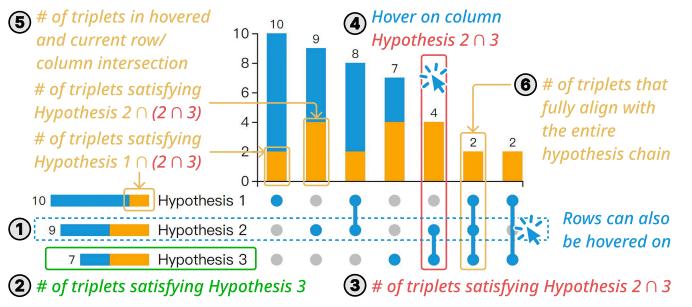


Fig. 3: UpSet Plot: Visualization of entity triplets alignment organized by intersections and inclusion relationships among hypotheses.

4.3 Pipeline Walkthrough

To clarify the pipeline’s structure, we walked through an SL prediction example together with **E1** and **E2**, systematically demonstrating how **Contextual Exploration**, **Hypothesis Generation**, and **Validation Selection** are applied in sequence, with detailed steps and explanations.

Contextual Exploration. At this stage (Fig. 4-①), the system provides experts with comprehensive and relevant information about their selected research focuses. It explains predictions with supplementary information on edge relationship granularity and integrates structured KG knowledge to support predictions. This structured information facilitates further mechanism exploration, potentially revealing the underlying mechanisms behind the predictions.

The process begins with training the GNN model on full-scale KG (Fig. 4-①), generating predictions along with interpretative paths. To assist experts in identifying research focuses matching their expertise and interest, while also filtering the most relevant predictions from vast results (**DR3**) (Fig. 4-②), the system incorporates a RAG-based retrieval framework with two modes: 1) Online LLM retrieval tailored to the biomedical domain. 2) RAG retrieval exclusively using local knowledge to prevent external data interference. Domain experts can engage with the RAG system (Fig. 4-⑧) by formulating queries based on their interests. For example, **E1**, an expert in breast cancer SL mechanisms, was interested in exploring SL in *salivary gland cancer* to identify potential mechanisms or patterns similar to *breast cancer*. To initiate this research, he queried RAG: “I would like to conduct research on *salivary gland cancer*, and I am an expert in *breast cancer*. Could you suggest some relevant predictions to facilitate the commencement of my study?” (Fig. 4-③) The system then returned **recommendations** for each research focus with the most relevant entities (Fig. 4-④), highlighted in the *Embedding View* (Fig. 4-⑤). Domain experts can refine their selection by integrating insights from RAG recommendations and the *Embedding View* using lasso selections (Fig. 4-⑥). As exemplified in the case, RAG first performed reasoning and local retrieval to generate a list of recommended genes. Among these, RAG recommended the gene *BRCA1* due to its association with both *salivary gland cancer* and *breast cancer*. Despite extensive focus on *BRCA1*, numerous unverified predictions persisted. Consequently, **E1** implemented a lasso selection on the cluster containing *BRCA1*, subsequently filtering *BRCA1* in the *Prediction View* (Fig. 4-⑦). Noticing many top-ranked predictions relied primarily on the *sl_gsg* relationship, which he considered biologically less meaningful, **E1** filtered out predictions whose interpretative paths consisted exclusively of *sl_gsg* connections. After that, he noticed that the prediction involving *BRCA1* and *USP1*, which was originally

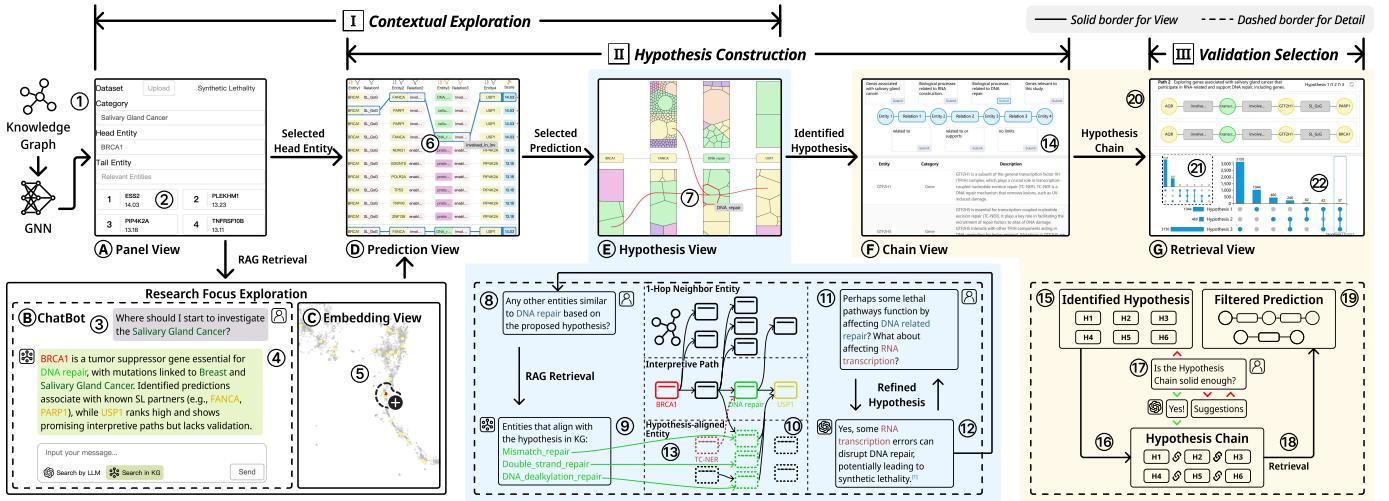


Fig. 4: The pipeline comprises three main components: **I Contextual Exploration**, **II Hypothesis Construction**, and **III Validation Selection**.

ranked 25th, had moved up to the 2nd position. Additionally, **E1** noticed experimentally validated genes *FANCA* and *PARP1* connected through two paths containing *DNA_repair*. Consequently, he identified *USP1* as a valuable prediction for further validation, completing research focus selection in the *Contextual Exploration* phase.

Once a research focus is chosen, experts can explore specific prediction paths (Fig. 4-⑥) in the *Hypothesis View* (Fig. 4-⑤). This view further displays selected paths along with their adjacent graph structures, offering contextual insights into the research topic (Fig. 4-⑦), aiding experts in comprehending the interpretative KG path and exploring additional related information. For example, **E1** queried for additional details regarding *DNA_repair* in the interpretative path (Fig. 4-⑧), obtaining a detailed explanation enhanced by external knowledge.

Then the LLM, with external information access, enriches interpretative paths by providing additional details in the *ChatBot*, allowing experts to bridge gaps in structured data by integrating common knowledge and cutting-edge insights not yet present in the KG. Experts can validate this information and, if valuable, integrate it into KG using text-to-KG methods within RAG, ensuring seamless KG expansion. Iterative querying and exploration allow experts to progressively refine their understanding of the predictions and the associated KG knowledge, laying the groundwork for subsequent hypothesis construction.

Hypothesis Construction. With prediction paths selected, domain experts require a solid basis to formulate hypotheses. Hence, the *Hypothesis View* leverages the LLM to generate an initial hypothesis regarding the underlying mechanisms. At this stage (Fig. 4-②), the *Hypothesis View* presents entities retrieved via RAG that align with the given hypothesis (Fig. 4-⑧, ⑨). Additionally, RAG paraphrases these descriptions, allowing domain experts to iteratively refine the reasoning until fully satisfied, ensuring an accurate understanding [36]. As demonstrated in the example, **E1** wished to explore further, so he shifted to the *Hypothesis View*. The LLM generated a potential hypothesis in the *ChatBot* that the target gene likely shares SL relationship with *DNA_repair*-related genes through interactions involving known SL genes. **E1** found this consistent with his domain knowledge and proposed that if *DNA_repair* could connect to different SL genes, similar entities might also yield similar predictions (Fig. 4-⑧). Consequently, the hypothesis was refined to encompass analogous entities related to *DNA_repair*. Then, the RAG retrieved three closely related entities: *Mismatch_repair*, *Double_strand_repair*, and *DNA_dealkylation_repair* (Fig. 4-⑨). **E1** then observed that these entities were also linked to many analogous predicted entities (Fig. 4-⑩), thus strengthening his belief in the hypothesis.

By leveraging these aligned entities, experts can concretize their hypothesis, examine related entities, and explore potential connection patterns predicted by the model (Fig. 4-⑪). This iterative process (**DR4**) deepens their insights, allowing them to refine and adjust hypotheses based on emerging findings (Fig. 4-⑪, ⑫). As the hypothesis

evolves, the displayed entities update accordingly.

Throughout the iterative refinement, the LLM using external information provides heuristic guidance (Fig. 4-⑫) in the *ChatBot*, helping domain experts assess hypothesis validity, identify potential inconsistencies, and enhance logical coherence. Once satisfied with a constructed hypothesis (Fig. 4-⑮), experts can integrate it into a *hypothesis chain* within the *Chain View* (Fig. 4-⑯) for further retrieval.

A *hypothesis chain* links multiple hypotheses, forming a structure similar to a triplet (Fig. 4-⑯). However, it extends beyond simple entity-relationship pairs by incorporating textual descriptions to facilitate communication and encoding of hypotheses. Experts can continually refine existing hypotheses or develop new ones with LLM-generated suggestions (Fig. 4-⑰). With each submission of hypothesis analysis, aligned predictions are marked by ★ in the *prediction view*, helping experts refine the chain. By constructing complex chains, they move beyond merely summarizing patterns or drawing simple analogies from model predictions. Instead, experts integrate expertise and insights into a flexible and targeted **Validation Selection** process, enabling deeper exploration of intricate underlying mechanisms.

As shown in the case, **E1** recalled that some RNA transcriptions are important for *DNA_repair*. To explore further, he queried LLM whether the proposed hypothesis held water (Fig. 4-⑪). The LLM confirmed a strong association between RNA and *DNA_repair* (Fig. 4-⑫). However, subsequent analysis using RAG revealed a particular relationship: *Transcription-Coupled Nucleotide-Excision-Repair* (TC-NER) (Fig. 4-⑬), a pathway where RNA polymerase triggers *DNA_repair* during transcription. Surprisingly, this relationship wasn't directly linked to *DNA_repair* or analogous entities in KG. This prompted **E1** to ask the LLM for clarification. The LLM confirmed this was indeed a scientifically valid connection. Further exploration showed this relationship was only associated with a specific subset of genes, devoid of any discernible connections to other entities. **E1** then incorporated the relationship into the KG through the text-to-KG integration within RAG and adjusted the hypothesis chain to (Fig. 4-⑭): [genes associated with *salivary gland cancer*] → [biological processes related to RNA construction] → [biological processes related to *DNA_repair*] → [genes relevant to this study] (Fig. 4-⑯), querying the LLM to ascertain the reasonableness of this chain (Fig. 4-⑰). After reasoning, the LLM recommended refining the relationship between RNA and *DNA_repair*, suggesting that the chain should be rephrased to indicate that RNA transcription is “related to” or “supports” *DNA_repair*-related entities, which is predicated on the observation that, in certain instances, the inactivation of specific genes actually stimulate RNA transcription, which might not result in an SL pair. **E1** considered this suggestion reasonable and refined the hypothesis chain accordingly.

Validation Selection. At this stage (Fig. 4-③), after domain experts have formulated coherent and well-reasoned hypotheses, these hypotheses are retrieved based on all the entities identified through

RAG to ensure retrieval accuracy and comprehensiveness (Fig. 4-18) and are then compared against predictions and the KG (Fig. 4-19, 20) within the *Retrieval View* (Fig. 4-C). This process helps determine whether similar conclusions have been previously validated or if relevant paths exist within the predictions.

Particularly, the retrieval process is guided by the entities proposed by RAG in the hypothesis chain (Fig. 4-19), with the retrieved predictions grouped based on their alignment with the hypothesis (Fig. 4-20). Domain experts can then examine these results to identify candidates for further experimental validation. Additionally, they can leverage the *UpSet Plot* (Fig. 1-G) to detect potential inconsistencies in the hypothesis during retrieval, gaining insights for further refinement and optimization. As reflected in the case, following the construction of the hypothesis chain, **E1** conducted a targeted retrieval (Fig. 4-18) based on the refined hypothesis chain. Analysis of the results revealed that while many predictions included RNA transcription entities as intermediates and linked to final outcomes via *DNA_repair*, none of them fully aligned with the intermediate hypothesis criteria (Fig. 4-21). Drawing from prior observations that certain RNA processes in the KG were predominantly connected to genes, **E1** expanded the hypothesis to incorporate gene-centric biological processes related to *DNA_repair*. Subsequent retrievals demonstrated strong alignment between predictions and the revised hypothesis chain (Fig. 4-22). Notably, the gene *EP300*—highly associated with *salivary gland cancer*—was connected to predicted entities *CYP2C9* and *RAD23B* through RNA-transcription-related entities *transcription-coupled_nucleotide-excision_repair* and *DNA_repair*-related genes *GTF2H1*, *GTF2H4*, and *GTF2H5*. This reinforced the hypothesis that *DNA_repair* mechanisms may exhibit SL correlations across disease contexts, with RNA transcription acting as a potential extension of these mechanisms. However, the current KG lacked robust representations of RNA transcriptions–*DNA_repair* relationships, underscoring the need for supplemental data integration. These findings revealed areas for improvement in the KG and provided a novel perspective for further research into mechanistic synergies.

5 EVALUATION

We conducted a case study and a user study to evaluate the effectiveness, workflow, and usability of *HypoChainer*.

5.1 Case Study: Drug Repurposing Exploration

E3 and **E4**, specialists in drug mechanisms and gene therapy, focus on repurposing antiepileptic drugs (e.g., those treating spasms) to address other diseases and evaluate novel therapies. Their goals are twofold: (1) identify new applications for existing drugs and (2) assess their potential in mitigating complex diseases. They followed the Co-discovery Learning Protocol [38], with one author guiding the session, **E3** operating the system, and **E4** discussing insights in real time.

To initiate this process, **E3** uploaded a drug repurposing dataset (Fig. 1-1) and queried the RAG module: “I am researching antiepileptic drug repurposing. Can you suggest potentially related diseases?” (Fig. 1-2). The system performed a multi-faceted analysis of disease relationships, identifying 5 high-relevance candidates, including *Huntington’s disease*, *Episodic ataxia type 5*, *Parkinson’s disease*, *Photosensitive tonic-clonic seizures*, and *Generalized tonic-clonic seizures*. These results were prioritized in both textual and *Embedding View*, with *Episodic ataxia type 5* flagged as a top recommendation (Fig. 1-3). While **E3** had prior expertise in *Photosensitive tonic-clonic seizures*, he used the *Embedding View*’s lasso tool to isolate the *Episodic ataxia* cluster (Fig. 1-4), opting to explore RAG’s novel suggestion. In the *Prediction View*, **E3** observed frequent top-ranked associations with the *CACNA1C* gene (Fig. 1-5). **E4** noted that this is a known regulator of voltage-gated calcium channels and a common target in epilepsy-related drug mechanisms. To deepen **E3**’s investigation, he selected the less familiar *Episodic ataxia type 5* prediction (Fig. 1-6), leveraging the system to bridge knowledge gaps and validate hypotheses.

Hypothesis Refinement Workflow. In the *Hypothesis View*, the *ChatBot* generated an initial mechanistic hypothesis via the LLM: [Modulating agent alters ion channel function] → [Impacts] → [Gene regulatory pathways (affected by calcium channels)] → [Drives] →

[Neuronal network dynamics] → [Manifests as] → [Episodic neural dysfunction] (Fig. 1-7). **E3** first validated this chain (Fig. 1-8) by confirming that retrieved entities aligned with the hypothesis and were supported by contextual KG evidence (Fig. 1-9). He observed that integrating existing KG edges revealed novel sub-paths (Fig. 1-10) and multiple connections between hypothesis-aligned entities and one-hop neighbors, reinforcing the chain’s coherence. Cross-referencing the *Prediction View*, he noted consistency between the hypothesis and most interpretative paths (Fig. 1-11). **E3** considered this as evidence that the hypothesis chain explains most repurposing predictions.

Critical Insight and Hypothesis Revision. However, **E3** noted a misalignment between the hypothesis and high-score predictions for *Huntington’s disease* (Fig. 1-12). By filtering and examining these predictions, he observed that the original hypothesis aligned only partially. Guided by the LLM’s explanations, **E4**’s further analysis of the remaining predictions revealed that some drugs were indicated for depression, while others were linked to neurodegenerative diseases. This finding led **E3** to infer that the model’s predictions regarding *Huntington’s disease* primarily reflect the symptomatic alleviation effects rather than an influence on or delay of its underlying etiology. To reconcile this, **E3** refined the hypothesis chain to: [Drugs treating motor dysfunction, depression, or neurodegeneration] → [related to] → [Interacting genes or pathways] → [Participated in] → [Processes related to motor dysfunction, depression, or neurodegeneration] → [Led to] → [Diseases associated with Huntington’s disease].

Hypothesis Challenge and Iterative Refinement. During hypothesis validation, **E3** noted that the revised chain obscured the majority of the predictions. However, in the *Hypothesis View*, he identified a link between *Huntington’s disease* and the *Trinucleotide Repeat Expansion* entity—a genetic mechanism associated with neurodegenerative disorders, as noted by **E4** (Fig. 5-1). Intrigued by its absence in the *Prediction View*, **E3** queried the LLM, which explained that *Trinucleotide Repeat Expansion* is a well-established genetic cause of *Huntington’s disease*, yet it did not surface in any *Huntington’s* predictions. Intrigued by this discrepancy, **E3** revised the hypothesis chain to prioritize predictions involving *Trinucleotide Repeat Expansion* and submitted the updated hypothesis chain via the for LLM validation. The LLM provided critical feedback: while *Trinucleotide Repeat Expansion* refers to abnormal DNA sequence elongations that disrupt gene expression or protein function, its absence in predictions likely reflects the lack of direct therapeutics. However, the LLM emphasized indirect associations in literature, linking the entity to broader processes like *DNA_repair* and *histone deacetylase (HDAC)* regulation.

Insight-Driven Revision. Guided by this feedback, **E3** refined the hypothesis to focus on *HDAC*-related therapeutic targets. A subsequent retrieval (Fig. 1-13) identified therapies modulating *HDAC* activity in diseases involving abnormal repeat expansions, such as *Parkinson’s disease* [34] and *Amyotrophic Lateral Sclerosis (ALS)* (Fig. 1-16). This alignment validated the revised chain (Fig. 1-14, 15).

Critical Discovery and Therapeutic Proposal. Re-examining *Huntington’s* predictions, **E3** uncovered a previously overlooked interpretative path: [*Entinostat* → [decreases activity of] → [*HDAC1* gene] → [interacts with] → [*Histone H4*] → [gene associated with condition] → [*Huntington’s disease*]]. This *HDAC*-associated interpretative path aligned with literature suggesting *HDAC*’s role in disease progression [22]. Synthesizing these insights, **E3** and **E4** proposed the design of a novel *cocktail therapy* combining three symptom-alleviating drugs with *HDAC* inhibitors to potentially delay *Huntington’s* progression. He further considered exploring *HDAC*-targeted gene therapies as a complementary avenue for further research.

Takeaway Message. In the case study, **E3** and **E4** iteratively revised

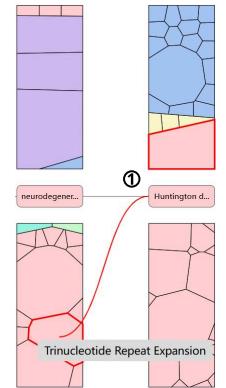


Fig. 5: ① **E3** noticed a link to *Huntington’s disease* from *Trinucleotide Repeat Expansion*, a 1-hop entity of *neurodegenerative disease*.

initial hypotheses as new insights emerged, reflecting a dynamic reasoning process. Although the final hypotheses diverged from the original, they represented logical extensions informed by system-driven exploration. For example, in repurposing drugs for *Huntington’s disease*, the expert initially focused on antiepileptic mechanisms for symptomatic relief. However, analysis of high-confidence predictions outside this scope revealed links to broader neurodegenerative pathways, prompting a shift toward a more integrative hypothesis. **E3** reported maintaining control over hypothesis development, with *HypoChainer* offering timely support when exploration stalled—clarifying complex predictions, suggesting new directions, and highlighting supporting KG evidence. **E4** praised the system’s alignment with conventional workflows, flexible entity categorization, and context-aware knowledge integration, which streamlined hypothesis refinement. They also noted that the system helped transcend initial cognitive frames and uncover unanticipated, yet scientifically valuable associations. This iterative refinement process underscores the system’s potential to support deeper and more structured hypothesis generation in complex discovery tasks.

5.2 User Study

We recruited 12 graduate students (Mean Age = 26.33, SD = 1.93; 6 males, 6 females) from bioinformatics or biomedical engineering backgrounds, including 6 PhD and 6 Master’s students. Participants were randomly assigned to either the Baseline or *HypoChainer* group, balanced by degree level (3 PhD and 3 Master’s per group).

Participants completed a drug repurposing task targeting *Hemophilia B*. To ensure objective evaluation, all were screened to confirm no prior knowledge of the disease or its mechanisms. Five system-generated interpretive paths were selected based on partial alignment with DrugMechDB-curated mechanisms of actions (MOAs) for *Eptacog Alfa* and *Nonacog Alfa*. The Baseline system (Appendix Fig. 8) used an LLM-only setup, omitting the *Hypothesis View*, to isolate the contributions of RAG and the view’s design while maintaining the overall workflow. Both drugs were excluded from training data to better simulate real-world discovery conditions. Participants received a 30-minute training session on task goals, system usage, and evaluation criteria, followed by 90 minutes to explore predictions, formulate hypotheses, and submit results. The process was screen-recorded. A task was deemed successful if participants retrieved ≤ 300 predictions and identified ≥ 3 of 5 reference predictions. A post-task questionnaire (Fig. 6) adapted from the *System Usability Scale* [8] assessed perceived system effectiveness, workflow, and usability.

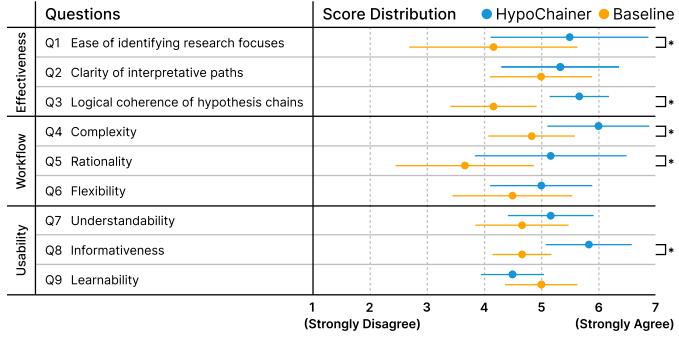


Fig. 6: The questionnaire results of the two systems in terms of system effectiveness, workflow and usability (* : $p < 0.05$).

In the Baseline group, 2 PhD participants (33.3%) completed the task within the time limit, compared to 4 in the *HypoChainer* group (3 PhDs, 1 Master’s; 66.7%). PhD participants generally outperformed Master’s students, likely due to stronger research backgrounds and faster system adaptation. *HypoChainer* showed clear advantages in *Effectiveness*, with significantly higher scores for *Ease of Identifying Research Directions* ($M = 5.50$, $SD = 1.38$, $p < 0.05$) and *Logical Coherence of Hypothesis Chains* ($M = 5.67$, $SD = 0.52$, $p < 0.05$) compared to the Baseline ($M = 4.17$, $SD = 1.47$; $M = 4.17$, $SD = 0.75$). Both groups scored similarly on *Clarity of Interpretative Paths* (*HypoChainer*: $M = 5.33$, $SD = 1.03$; *Baseline*: $M = 5.00$, $SD = 0.89$). For *Workflow*, *HypoChainer* scored higher on *Rationality* ($M = 5.17$, $SD = 1.33$,

$p < 0.05$) and *Flexibility* ($M = 5.00$, $SD = 0.89$, $p < 0.05$), supported by user feedback indicating fewer hallucinations. However, its integrated features—transitioning among LLM, RAG, and the *Hypothesis View*—were perceived as more complex ($M = 6.00$, $SD = 0.89$) than the Baseline ($M = 4.83$, $SD = 0.75$). In *Usability*, *HypoChainer* had a slightly lower score in *Learnability* ($M = 4.50$, $SD = 0.55$) than the Baseline ($M = 5.00$, $SD = 0.63$), suggesting a steeper learning curve. Nevertheless, it received significantly higher ratings for *Informativeness* ($M = 5.83$, $SD = 0.75$, $p < 0.05$) and *Understandability* ($M = 5.17$, $SD = 0.75$, $p < 0.05$), due to its clear reasoning structure and entity-level explanations. Overall, *HypoChainer* provided more effective, transparent, and insightful support for hypothesis generation, despite a slightly higher learning threshold.

6 DISCUSSION AND LIMITATION

Lessons Learned. During the evaluation of RAGs, we identified critical trade-offs between accuracy and computational cost. Tests revealed that response accuracy depends on the scale of the local KG and the entity retrieval limit per query. Increasing the retrieval limit from the default to more entities markedly improved accuracy but incurred higher token consumption and prolonged query times. This highlights the necessity of balancing query performance and cost in large-scale RAG-based workflows. Additionally, while model predictions occasionally diverged from the constructed hypothesis chains, domain experts emphasized that such discrepancies do not diminish the validity of the chains themselves. They highlighted that rigorously derived hypotheses—even without full alignment with retrieval results—retain significant values, particularly in uncovering novel insights.

Generalization and Scalability. The hypothesis-driven architecture has demonstrated broad applicability across diverse research domains, particularly in scientific discovery tasks involving KG-based node and link prediction. When integrated with RAG, the system offers a cost-effective and flexible alternative to domain-specific fine-tuning, enabling more efficient adaptation to new research areas. Given the generalizability of RAG and the system’s modular KG-based design, most components are transferable across domains, requiring only the substitution of domain-specific KGs and predictive models. Domain experts have also identified promising applications in areas such as protein structure and function prediction via interaction networks, as well as novel materials discovery. However, the system’s performance is largely dependent on the quality of the underlying KGs and the effectiveness of the associated predictive models. As advances in text-based KG extraction continue to improve accuracy and robustness, the system’s cross-domain applicability is expected to expand further.

Limitations. While RAG supports text-to-KG conversion, the process remains resource-intensive and prone to inaccuracies, especially in data-rich domains. To reduce risk, we restricted integration to small, verified updates using public KGs. However, the slow pace of KG updates compared to rapid scientific progress highlights the need for more accurate, automated extraction methods. Although RAG helps mitigate hallucinations by retrieving from traceable sources, such issues persist. Future improvements could include fact-verification modules—e.g., cross-referencing authoritative sources or applying confidence thresholds to flag uncertain outputs. Given the complexity of hypothesis-driven discovery, some interaction complexity is unavoidable, though we anticipate continued simplification as AI reliability advances.

7 CONCLUSION AND FUTURE WORK

This study introduces *HypoChainer*, a collaborative framework that synergizes LLMs and KGs to advance hypothesis-driven scientific discovery. *HypoChainer* provides three main functionalities: **Contextual Exploration**, **Hypothesis Construction** and **Validation Selection**. A case study and expert interviews demonstrate *HypoChainer*’s capability to synthesize context-aware knowledge efficiently, construct and refine hypothesis chains systematically, and facilitate informed validation selections. Future work includes enhancing text-KG integration through more robust methods and accelerating knowledge discovery via more workflow automation, while preserving critical human oversight to ensure scientific rigor and actionable insights.

ACKNOWLEDGMENTS

We gratefully acknowledge Dr. Jia Liu and Dr. Yifeng Yang from the Institute of Immunochemistry at ShanghaiTech University for their valuable collaboration, as well as the anonymous reviewers for their insightful feedback. This research was supported by the National Natural Science Foundation of China (No. 62372298), the “AI Technologies for Accelerating Biopharmaceutical R&D – School of Information Science and Technology” (No. 2024X0203-902-01), the Shanghai Engineering Research Center of Intelligent Vision and Imaging, the Shanghai Frontiers Science Center of Human-centered Artificial Intelligence (ShangHAI), and the MoE Key Laboratory of Intelligent Perception and Human-Machine Collaboration (KLIP-HuMaCo). Part of the experimental work was conducted with support from the Core Facility Platform of Computer Science and Communication, SIST, ShanghaiTech University.

REFERENCES

- [1] A. Ahmed, M. Saleem, M. Alzeen, et al. Leveraging large language models to enhance machine learning interpretability and predictive performance: A case study on emergency department returns for mental health patients, 2025. doi: <https://doi.org/10.48550/arXiv.2502.00025> 2
- [2] S. Amer-Yahia, A. Bonifati, L. Chen, G. Li, X. Shim, Kyuseok, et al. From large language models to databases and back: A discussion on research and education. *SIGMOD Rec.*, 52(3):49–56, 8 pages, Nov. 2023. doi: [10.1145/3631504.3631518](https://doi.org/10.1145/3631504.3631518) 3
- [3] L. Asprino, C. Colonna, M. Mongiovì, M. Porena, and V. Presutti. Pattern-based visualization of knowledge graphs, 2021. doi: [10.48550/arXiv.2106.12857](https://doi.org/10.48550/arXiv.2106.12857) 3
- [4] J. Baek, S. K. Jauhar, S. Cucerzan, and S. J. Hwang. ResearchAgent: Iterative research idea generation over scientific literature with large language models, 2025. doi: [10.48550/arXiv.2404.07738](https://doi.org/10.48550/arXiv.2404.07738) 3
- [5] M. Balzer and O. Deussen. Voronoi treemaps. In *IEEE Symposium on Information Visualization, 2005. INFOVIS 2005.*, pp. 49–56, 2005. doi: [10.1109/INFVIS.2005.1532128](https://doi.org/10.1109/INFVIS.2005.1532128) 3, 6
- [6] M. R. Birtwistle. Analytical reduction of combinatorial complexity arising from multiple protein modification sites. *Journal of The Royal Society Interface*, 12(103):20141215, 2015. doi: [10.1098/rsif.2014.1215](https://doi.org/10.1098/rsif.2014.1215) 2
- [7] T. Boger, S. B. Most, and S. L. Franconeri. Jurassic mark: Inattentional blindness for a datasaurus reveals that visualizations are explored, not seen. In *2021 IEEE Visualization Conference (VIS)*, pp. 71–75, 2021. doi: [10.1109/VIS49827.2021.9623273](https://doi.org/10.1109/VIS49827.2021.9623273) 2
- [8] J. Brooke. SUS: a retrospective. *Journal of Usability Studies*, 8:29–40, 01 2013. 9
- [9] Y. Chang, X. Wang, J. Wang, Y. Wu, L. Yang, K. Zhu, et al. A survey on evaluation of large language models. *ACM transactions on intelligent systems and technology*, 15(3):1–45, 2024. doi: [10.1145/3641289](https://doi.org/10.1145/3641289) 2
- [10] S. Chen, S. Li, S. Chen, and X. Yuan. R-Map: A map metaphor for visualizing information reposting process in social media. *IEEE Transactions on Visualization and Computer Graphics*, 26(1):1204–1214, 2020. doi: [10.1109/TVCG.2019.2934263](https://doi.org/10.1109/TVCG.2019.2934263) 3
- [11] N. Choudhary and C. K. Reddy. Complex logical reasoning over knowledge graphs using large language models, 2024. doi: [10.48550/arXiv.2305.01157](https://doi.org/10.48550/arXiv.2305.01157) 3
- [12] V. Clarke and V. Braun. Thematic analysis. *The journal of positive psychology*, 12:297–298, 2017. doi: [10.1080/17439760.2016.1262613](https://doi.org/10.1080/17439760.2016.1262613) 3
- [13] T. Dang, P. Murray, and A. Forbes. BioLinker: Bottom-up exploration of protein interaction networks. In *2017 IEEE Pacific Visualization Symposium (PacificVis)*, pp. 265–269. IEEE, 2017. doi: [10.1109/PACIFICVIS.2017.8031603](https://doi.org/10.1109/PACIFICVIS.2017.8031603) 2, 3
- [14] S. Deng, S. Wang, H. Rangwala, L. Wang, and Y. Ning. Cola-GNN: Cross-location Attention based Graph Neural Networks for Long-term ILI Prediction. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, CIKM ’20, 10 pages, p. 245–254, 2020. doi: [10.1145/3340531.3411975](https://doi.org/10.1145/3340531.3411975) 2
- [15] D. Edge, H. Trinh, N. Cheng, et al. From Local to Global: A graph RAG approach to query-focused summarization, 2025. doi: [10.48550/arxiv.2404.16130](https://doi.org/10.48550/arxiv.2404.16130) 5
- [16] A. Ghafarollahi and M. J. Buehler. SciAgents: Automating scientific discovery through bioinspired multi-agent intelligent graph reasoning. *Advanced Materials*, p. 2413523. doi: [10.1002/adma.202413523](https://doi.org/10.1002/adma.202413523) 2, 3
- [17] J. Gottweis, W.-H. Weng, A. Daryin, et al. Towards an AI co-scientist, 2025. doi: [10.48550/arXiv.2502.18864](https://doi.org/10.48550/arXiv.2502.18864) 3
- [18] S. Gratzl, A. Lex, N. Gehlenborg, H. Pfister, and M. Streit. LineUp: Visual analysis of multi-attribute rankings. *IEEE Transactions on Visualization and Computer Graphics*, 19(12):2277–2286, 2013. doi: [10.1109/TVCG.2013.173](https://doi.org/10.1109/TVCG.2013.173) 5
- [19] Z. Guo, L. Xia, Y. Yu, T. Ao, and C. Huang. LightRAG: Simple and fast retrieval-augmented generation, 2024. doi: [10.48550/arXiv.2410.05779](https://doi.org/10.48550/arXiv.2410.05779) 5
- [20] Y. Hasin, M. Seldin, and A. Lusis. Multi-omics approaches to disease. *Genome biology*, 18:1–15, 2017. doi: [10.1186/s13059-017-1215-1](https://doi.org/10.1186/s13059-017-1215-1) 3
- [21] L. He, S. Zheng, T. Yang, and F. Zhang. KLMo: Knowledge graph enhanced pretrained language model with fine-grained relationships. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pp. 4536–4542. Association for Computational Linguistics, Nov. 2021. doi: [10.18653/v1/2021.findings-emnlp.384](https://doi.org/10.18653/v1/2021.findings-emnlp.384) 3
- [22] K. Hecklau, S. Mueller, S. P. Koch, et al. The effects of selective inhibition of histone deacetylase 1 and 3 in Huntington’s disease mice. *Frontiers in molecular neuroscience*, 14:616886, 2021. doi: [10.3389/fnmol.2021.616886](https://doi.org/10.3389/fnmol.2021.616886) 8
- [23] S. Henry and B. T. McInnes. Literature based discovery: models, methods, and trends. *Journal of biomedical informatics*, 74:20–32, 2017. doi: [10.1016/j.jbi.2017.08.011](https://doi.org/10.1016/j.jbi.2017.08.011) 2
- [24] X. Hu, H. Fu, J. Wang, Y. Wang, Z. Li, R. Xu, Y. Lu, Y. Jin, L. Pan, and Z. Lan. Nova: An iterative planning and search approach to enhance novelty and diversity of LLM generated ideas, 2024. doi: [10.48550/arXiv.2410.14255](https://doi.org/10.48550/arXiv.2410.14255) 3
- [25] D. Huang, C. Yan, Q. Li, and X. Peng. From large language models to large multimodal models: A literature review. *Applied Sciences*, 14(12), 2024. doi: [10.3390/app14125068](https://doi.org/10.3390/app14125068) 2
- [26] G. Izocard and E. Grave. Leveraging passage retrieval with generative models for open domain question answering. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pp. 874–880. Association for Computational Linguistics, Apr. 2021. doi: [10.18653/v1/2021.eacl-main.74](https://doi.org/10.18653/v1/2021.eacl-main.74) 3
- [27] H. Jiang, S. Shi, S. Zhang, J. Zheng, and Q. Li. SLInterpreter: An exploratory and iterative Human-AI collaborative system for GNN-based synthetic lethal prediction. *IEEE Transactions on Visualization and Computer Graphics*, 31(1):919–929, 2025. doi: [10.1109/TVCG.2024.3456325](https://doi.org/10.1109/TVCG.2024.3456325) 2, 3, 6
- [28] Z. Jiang, F. F. Xu, J. Araki, and G. Neubig. How can we know what language models know? *Transactions of the Association for Computational Linguistics*, 8:423–438, 2020. doi: [10.1162/tacl_a_00324](https://doi.org/10.1162/tacl_a_00324) 3
- [29] Z. Jin, Y. Wang, Q. Wang, Y. Ming, T. Ma, and H. Qu. GNNLens: A visual analytics approach for prediction error diagnosis of graph neural networks. *IEEE Transactions on Visualization and Computer Graphics*, 29(6):3024–3038, 2023. doi: [10.1109/TVCG.2022.3148107](https://doi.org/10.1109/TVCG.2022.3148107) 3
- [30] T. Kaufmann, P. Weng, V. Bengs, and E. Hüllermeier. A survey of reinforcement learning from human feedback, 2024. doi: [10.48550/arXiv.2312.14925](https://doi.org/10.48550/arXiv.2312.14925) 3
- [31] J. Kehrer, F. Ladstädter, P. Muigg, H. Doleisch, A. Steiner, and H. Hauser. Hypothesis generation in climate research with interactive visual data exploration. *IEEE Transactions on Visualization and Computer Graphics*, 14(6):1579–1586, 2008. doi: [10.1109/TVCG.2008.139](https://doi.org/10.1109/TVCG.2008.139) 2
- [32] M. Krenn, L. Buffoni, B. Coutinho, et al. Forecasting the future of artificial intelligence with machine learning-based link prediction in an exponentially growing knowledge network. *Nature Machine Intelligence*, 5(11):1326–1335, 2023. doi: [10.1038/s42256-023-00735-0](https://doi.org/10.1038/s42256-023-00735-0) 2
- [33] S. Kumar, T. Ghosal, V. Goyal, and A. Ekbal. Can large language models unlock novel scientific research ideas?, 2024. doi: [10.48550/arXiv.2409.06185](https://doi.org/10.48550/arXiv.2409.06185) 2
- [34] V. Kumar. Understanding the role of histone deacetylase and their inhibitors in neurodegenerative disorders: Current targets and future perspective. *Current Neuropharmacology*, 20(1):158–178, 2022. doi: [10.2174/1570159X19666210609160017](https://doi.org/10.2174/1570159X19666210609160017) 8
- [35] A. Lex, N. Gehlenborg, H. Strobelt, R. Vuillemot, and H. Pfister. Upset: Visualization of intersecting sets. *IEEE Transactions on Visualization and Computer Graphics*, 20(12):1983–1992, 2014. doi: [10.1109/TVCG.2014.2346248](https://doi.org/10.1109/TVCG.2014.2346248) 6
- [36] H. Li, G. Appleby, and A. Suh. Linkq: An LLM-assisted visual interface for knowledge graph question-answering. In *2024 IEEE Visualization and Visual Analytics (VIS)*, pp. 116–120, 2024. doi: [10.1109/VIS55277.2024.000317](https://doi.org/10.1109/VIS55277.2024.000317) 7
- [37] L. Liang, M. Sun, Z. Gui, et al. KAG: Boosting LLMs in professional domains via knowledge augmented generation, 2024. doi: [10.48550/arXiv.2409.13731](https://doi.org/10.48550/arXiv.2409.13731) 3

- [38] K. H. Lim, L. M. Ward, and I. Benbasat. An empirical study of computer system learning: Comparison of co-discovery and self-discovery methods. *Information Systems Research*, 8(3):254–272, 1997. doi: [10.1287/isre.8.3.254](https://doi.org/10.1287/isre.8.3.254) 8
- [39] X. Lin, Z. Quan, Z.-J. Wang, T. Ma, and X. Zeng. KGNN: Knowledge graph neural network for drug-drug interaction prediction. In C. Bessiere, ed., *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20*, pp. 2739–2745. International Joint Conferences on Artificial Intelligence Organization, 7 2020. Main track. doi: [10.24963/ijcai.2020/380](https://doi.org/10.24963/ijcai.2020/380) 2
- [40] C. Lu, C. Lu, R. T. Lange, J. Foerster, J. Clune, and D. Ha. The AI scientist: Towards fully automated open-ended scientific discovery, 2024. doi: [10.48550/arXiv.2408.06292](https://doi.org/10.48550/arXiv.2408.06292) 3
- [41] C. Ma, Z. Zhou, H. Liu, and D. Koslicki. KGML-xDTD: a knowledge graph-based machine learning framework for drug treatment prediction and mechanism description. *GigaScience*, 12:giad057, 08 2023. doi: [10.1093/gigascience/giad057](https://doi.org/10.1093/gigascience/giad057) 2, 5
- [42] P. Ma, T.-H. Wang, M. Guo, Z. Sun, J. B. Tenenbaum, D. Rus, C. Gan, and W. Matusik. LLM and simulation as bilevel optimizers: A new paradigm to advance physical scientific discovery, 2024. doi: [10.48550/arXiv.2405.09783](https://doi.org/10.48550/arXiv.2405.09783) 3
- [43] L. McInnes and J. Healy. UMAP: Uniform manifold approximation and projection for dimension reduction. *ArXiv*, abs/1802.03426, 2018. doi: [10.21105/joss.00861](https://doi.org/10.21105/joss.00861) 5
- [44] H. Naveed, A. U. Khan, S. Qiu, M. Saqib, S. Anwar, M. Usman, N. Akhtar, N. Barnes, and A. Mian. A comprehensive overview of large language models, 2024. doi: [10.48550/arXiv.2307.06435](https://doi.org/10.48550/arXiv.2307.06435) 2
- [45] H. Nigam, M. Patwardhan, L. Vig, and G. Shroff. Acceleron: A tool to accelerate research ideation, 2024. doi: [10.48550/arXiv.2403.04382](https://doi.org/10.48550/arXiv.2403.04382) 3
- [46] S. Paley, R. Billington, J. Herson, M. Krummenacker, and P. D. Karp. Pathway tools visualization of organism-scale metabolic networks. *Metabolites*, 11(2):64, 2021. doi: [10.3390/metabolite11020064](https://doi.org/10.3390/metabolite11020064) 3
- [47] S. Pan, L. Luo, Y. Wang, C. Chen, J. Wang, and X. Wu. Unifying large language models and knowledge graphs: A roadmap. *IEEE Transactions on Knowledge and Data Engineering*, 2024. doi: [10.48550/arXiv.2306.08302](https://doi.org/10.48550/arXiv.2306.08302) 3
- [48] C. Peng, F. Xia, M. Naseriparsa, and F. Osborne. Knowledge graphs: Opportunities and challenges. *Artificial Intelligence Review*, 56(11):13071–13102, November 2023. doi: [10.1007/s10462-023-10465-9](https://doi.org/10.1007/s10462-023-10465-9) 2
- [49] J. Pérez, M. Arenas, and C. Gutierrez. Semantics and complexity of sparql. *ACM Trans. Database Syst.*, 34(3), article no. 16, 45 pages, Sept. 2009. doi: [10.1145/1567274.1567278](https://doi.org/10.1145/1567274.1567278) 3
- [50] G. Perković, A. Drobniak, and I. Botički. Hallucinations in LLMs: Understanding and addressing challenges. In *2024 47th MIPRO ICT and Electronics Convention (MIPRO)*, pp. 2084–2088, 2024. doi: [10.1109/MIPRO60963.2024.10569238](https://doi.org/10.1109/MIPRO60963.2024.10569238) 2, 3
- [51] K. Popper. *The logic of scientific discovery*. Routledge, 2005. doi: [10.4324/9780203994627](https://doi.org/10.4324/9780203994627) 2
- [52] K. Pu, K. J. K. Feng, T. Grossman, T. Hope, B. D. Mishra, M. Latzke, J. Bragg, J. C. Chang, and P. Siangliulue. IdeaSynth: Iterative research idea development through evolving and composing idea facets with literature-grounded feedback, 2024. doi: [10.48550/arXiv.2410.04025](https://doi.org/10.48550/arXiv.2410.04025) 2
- [53] B. Qi, K. Zhang, H. Li, K. Tian, S. Zeng, Z.-R. Chen, and B. Zhou. Large language models are zero shot hypothesis proposers, 2023. doi: [10.48550/arXiv.2311.05965](https://doi.org/10.48550/arXiv.2311.05965) 2
- [54] W. Qiang and Z. Zhongli. Reinforcement learning model, algorithms and its application. In *2011 International Conference on Mechatronic Science, Electric Engineering and Computer (MEC)*, pp. 1143–1146, 2011. doi: [10.1109/MEC.2011.6025669](https://doi.org/10.1109/MEC.2011.6025669) 3
- [55] T. Schimanski, J. Ni, M. Kraus, E. Ash, and M. Leippold. Towards faithful and robust LLM specialists for evidence-based question-answering. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1913–1931. Association for Computational Linguistics, Aug. 2024. doi: [10.18653/v1/2024.acl-long.105](https://doi.org/10.18653/v1/2024.acl-long.105) 2, 4
- [56] P. Schneider, N. Rehtanz, K. Jokinen, and F. Matthes. From data to dialogue: Leveraging the structure of knowledge graphs for conversational exploratory search, 2023. doi: [10.48550/arXiv.2310.05150](https://doi.org/10.48550/arXiv.2310.05150) 3
- [57] C. Shi, F. Nie, Y. Hu, Y. Xu, L. Chen, X. Ma, and Q. Luo. MedChemLens: An interactive visual tool to support direction selection in interdisciplinary experimental research of medicinal chemistry. *IEEE Transactions on Visualization and Computer Graphics*, 29(1):63–73, 2022. doi: [10.1109/TVCG.2022.3209434](https://doi.org/10.1109/TVCG.2022.3209434) 2, 4
- [58] C. Si, D. Yang, and T. Hashimoto. Can LLMs generate novel research ideas? a large-scale human study with 100+ NLP researchers, 2024. doi: [10.48550/arXiv.2409.04109](https://doi.org/10.48550/arXiv.2409.04109) 2
- [59] H. Su, R. Chen, S. Tang, X. Zheng, J. Li, Z. Yin, W. Ouyang, and N. Dong. Two heads are better than one: A multi-agent system has the potential to improve scientific idea generation. *IEEE Transactions on Cognitive and Developmental Systems*, 2024. doi: [10.1109/TCDS.2025.3530945](https://doi.org/10.1109/TCDS.2025.3530945) 3
- [60] K. M. Tolle, D. S. W. Tansley, and A. J. G. Hey. *The Fourth Paradigm: Data-intensive Scientific Discovery*. 2009. doi: [10.1109/JPROC.2011.2155130](https://doi.org/10.1109/JPROC.2011.2155130) 2
- [61] D. R. Unni, S. A. T. Moxon, M. Bada, et al. Biolink model: A universal schema for knowledge graphs in clinical, biomedical, and translational science. *Clinical and Translational Science*, 15(8):1848–1855, 2022. doi: [10.1111/cts.13302](https://doi.org/10.1111/cts.13302) 5
- [62] S. K. Vohra, P. Harth, Y. Isoe, A. Bahl, H. Fotowat, F. Engert, H.-C. Hege, and D. Baum. A visual interface for exploring hypotheses about neural circuits. *IEEE Transactions on Visualization and Computer Graphics*, 30(7):3945–3958, 2024. doi: [10.1109/TVCG.2023.3243668](https://doi.org/10.1109/TVCG.2023.3243668) 2
- [63] D. Vrandečić and M. Krötzsch. Wikidata: a free collaborative knowledgebase. *Communications of the ACM*, 57(10):78–85, 2014. doi: [10.1145/2629489](https://doi.org/10.1145/2629489) 3
- [64] F. Wang, X. Zhou, W. Hu, Z. Luo, W. Luo, and X. Bai. LLM assists hypothesis generation and testing for deliberative questions. 13 pages, p. 424–436. Springer-Verlag, 2024. doi: [10.1007/978-981-97-9434-8_33](https://doi.org/10.1007/978-981-97-9434-8_33) 2
- [65] Q. Wang, K. Huang, P. Chanda, M. Zitnik, and N. Gehlenborg. Extending the nested model for user-centric XAI: A design study on GNN-based drug repurposing. *IEEE Transactions on Visualization and Computer Graphics*, 29(1):1266–1276, 2023. doi: [10.1109/TVCG.2022.3209435](https://doi.org/10.1109/TVCG.2022.3209435) 3
- [66] T. Wang, S. Chen, Y. Wang, et al. From in Silico to in Vitro: A comprehensive guide to validating bioinformatics findings, 2025. doi: [10.48550/arXiv.2410.07076](https://doi.org/10.48550/arXiv.2410.07076) 2, 3
- [67] X. Wang, T. Gao, Z. Zhu, Z. Zhang, Z. Liu, J. Li, and J. Tang. KEPLER: A unified model for knowledge embedding and pre-trained language representation. *Transactions of the Association for Computational Linguistics*, 9:176–194, 2021. doi: [10.1162/tacl_a_00360](https://doi.org/10.1162/tacl_a_00360) 3
- [68] Y. Yan, Y. Hou, Y. Xiao, R. Zhang, and Q. Wang. KNOWNET: Guided health information seeking from LLMs via knowledge graph integration. *IEEE Transactions on Visualization and Computer Graphics*, 2024. doi: [10.1109/TVCG.2024.3456364](https://doi.org/10.1109/TVCG.2024.3456364) 3
- [69] Y. Yang, Y. Wang, Y. Li, S. Sen, L. Li, and Q. Liu. Unleashing the potential of large language models for predictive tabular tasks in data science. 2024. doi: [10.18653/v1/2024.findings-emnlp.224](https://doi.org/10.18653/v1/2024.findings-emnlp.224) 2
- [70] R. Yao, Z. Shen, X. Xu, G. Ling, R. Xiang, T. Song, F. Zhai, and Y. Zhai. Knowledge mapping of graph neural networks for drug discovery: a bibliometric and visualized analysis. *Frontiers in Pharmacology*, 15, 2024. doi: [10.3389/fphar.2024.1393415](https://doi.org/10.3389/fphar.2024.1393415) 2
- [71] J. Yuan, X. Yan, B. Shi, T. Chen, W. Ouyang, B. Zhang, et al. Dolphin: Closed-loop open-ended auto-research through thinking, practice, and feedback, 2025. doi: [10.48550/arXiv.2501.03916](https://doi.org/10.48550/arXiv.2501.03916) 3
- [72] K. Zhang, M. Wu, Y. Liu, Y. Feng, and J. Zheng. KR4SL: knowledge graph reasoning for explainable prediction of synthetic lethality. *Bioinformatics*, 39:i158–i167, 06 2023. doi: [10.1093/bioinformatics/btad261](https://doi.org/10.1093/bioinformatics/btad261) 2, 5
- [73] Y. Zhang, B. Hu, Z. Chen, L. Guo, Z. Liu, Z. Zhang, L. Liang, H. Chen, and W. Zhang. Multi-domain knowledge graph collaborative pre-training and prompt tuning for diverse downstream tasks, 2024. doi: [10.48550/arXiv.2405.13085](https://doi.org/10.48550/arXiv.2405.13085) 2
- [74] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong, Y. Du, C. Yang, Y. Chen, Z. Chen, et al. A survey of large language models, 2025. doi: [10.48550/arXiv.2303.18223](https://doi.org/10.48550/arXiv.2303.18223) 2
- [75] C. Zheng, Y. Zhang, Z. Huang, C. Shi, M. Xu, and X. Ma. Disciplink: Unfolding Interdisciplinary Information Seeking Process via Human-AI Co-Exploration. In *Proceedings of the 37th Annual ACM Symposium on User Interface Software and Technology*, pp. 1–20, 2024. doi: [10.1145/3654777.3676366](https://doi.org/10.1145/3654777.3676366) 3
- [76] T. Zheng, Z. Deng, H. T. Tsang, W. Wang, J. Bai, Z. Wang, and Y. Song. From automation to autonomy: A survey on large language models in scientific discovery, 2025. 2
- [77] Y. Zhou, H. Liu, T. Srivastava, H. Mei, and C. Tan. Hypothesis generation with large language models. 2024. doi: [10.18653/v1/2024.nlp4science-1.10](https://doi.org/10.18653/v1/2024.nlp4science-1.10) 2, 3
- [78] Z. Zhou, X. Feng, L. Huang, et al. From hypothesis to publication: A comprehensive survey of AI-driven research support systems, 2025. doi: [10.48550/arXiv.2310.05150](https://doi.org/10.48550/arXiv.2310.05150) 2