# Scientific Hypothesis Generation and Validation: Methods, Datasets, and Future Directions

Adithya Kulkarni[1], Fatimah Alotaibi[1], Xinyue Zeng[1], Longfeng Wu[1], Tong Zeng[1], Barry Menglong Yao[1], Minqian Liu[1], Shuaicheng Zhang[1], Dawei Zhou[1], and Lifu Huang[2]

[1]Virginia Tech, Blacksburg, VA, USA
[2]University of California, Davis, CA, USA

## 1 Introduction

Scientific discovery has long been the subject of computational modeling, particularly through systems that frame discovery as a structured problem-solving process grounded in cognitive science and artificial intelligence Bradshaw et al. [1983], Simon [1992], Langley and Jones [1988], Langley [1998, 2000], Džeroski et al., Langley and Simon [2013], Langley [2024], Bengio and LeCun [2007], Hinton et al. [2006], Boden [1998]. These early approaches emphasized the iterative formulation and refinement of hypotheses using heuristic rules, domain knowledge, and symbolic representations. While such systems successfully simulated elements of scientific reasoning, ranging from rediscovering physical laws to inferring causal structures, their scalability and adaptability to unstructured data remained limited. In parallel, the advent of Large Language Models (LLMs) AI@Meta [2024], Hurst et al. [2024] has marked a transformative shift in scientific knowledge creation and validation. These models, trained on expansive corpora encompassing text, numerical data, and multimodal inputs, possess the remarkable capability to synthesize diverse datasets, identify latent patterns, and accelerate the hypothesis generation process at an unprecedented scale Beltagy et al. [2019], Qi et al. [2024], Wang et al. [2023, 2025], Ding et al. [2025], Ji et al. [2024], Goodfellow et al. [2016]. Moreover, LLMs have demonstrated proficiency in extracting meaningful relationships from unstructured text, facilitating hypothesis discovery in fields such as biomedical research and materials science Sybrandt et al. [2018], Ghafarollahi and Buehler [2024a], Lin et al. [2025], Bran et al. [2023], Kadic et al. [2019], Bertoldi et al. [2017], Jia et al. [2020], Jiao and Alavi [2021], Bauer et al. [2017], Papadimitriou et al. [2024]. These capabilities make LLMs instrumental in confronting the complexity and scale of contemporary scientific inquiry, enabling a paradigm where data-driven insights complement and extend human reasoning in the discovery process Yang et al. [2023], Shojaee et al. [2024], Romera-Paredes et al. [2024], Trinh et al. [2024], Lesica et al. [2021].

The growing capabilities of LLMs underscore their potential to revolutionize hypothesis generation and validation. Frameworks such as The AI Scientist Lu et al. [2024] and SciAgents Ghafarollahi and Buehler [2024b] epitomize the advancements in agentic AI systems that autonomously undertake significant elements of the scientific process, including experimental validation and the drafting of manuscripts. These systems harness advanced methodologies such as retrieval-augmented generation (RAG) Chen et al. [2024], knowledge graph integration Zhou et al. [2024], Sybrandt et al. [2017], and causal inference Jha et al. [2019], enabling the generation of hypotheses that are not only testable but also interdisciplinary. By systematically mapping connections across seemingly unrelated domains, LLMs uncover insights that human researchers might overlook due to cognitive constraints or disciplinary silos Fawzi et al. [2022], Touvron et al. [2023], Qi et al. [2023]. This integration of machine learning models and LLMs into the hypothesis generation process redefines the boundaries of scientific exploration, opening avenues for cross-domain innovations that were

---

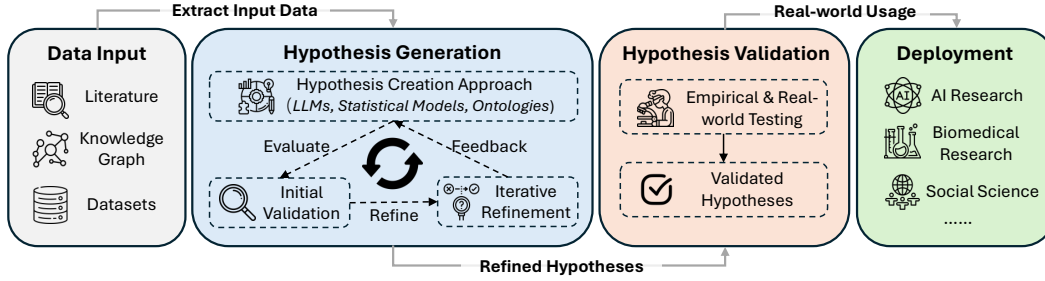*Correspondence to: `aditkulk@vt.edu`, `lfuhuang@ucdavis.edu`

Figure 1: Overview of the scientific hypothesis generation and validation pipeline integrating LLMs, statistical models, and ontologies. The figure illustrates the stages from data input and hypothesis creation to iterative validation and real-world deployment, highlighting the feedback loops that refine hypotheses over time.

previously inconceivable. Figure 1 provides an overview of the scientific hypothesis generation and validation pipeline.

Table 1: Examples of LLM-Driven Tools for Hypothesis Creation and Validation

| Tool | Application | Contribution to Scientific Discovery |
|---|---|---|
| *AlphaFold* Jumper et al. [2021] | Protein structure prediction | Resolved a decades-long challenge in biology, accelerating drug discovery |
| *Crispr-GPT* Huang et al. [2024] | Gene-editing experiment design | Automated hypothesis creation for CRISPR-based research |
| *SciAgents* Ghafarollahi and Buehler [2024b] | Dynamic knowledge graph generation | Maps relationships between interdisciplinary concepts, revealing unexplored connections |
| *Discovering Faster Matrix Multiplication Algorithms* Fawzi et al. [2022] | Algorithmic optimization | Demonstrates how reinforcement learning can refine complex mathematical hypotheses |
| *Materials Project* Jain et al. [2020] | Materials property prediction | Enables hypothesis generation about novel material structures and properties |
| *MOLIERE* Sybrandt et al. [2017, 2018] | Biomedical hypothesis validation | Retrospectively tests hypotheses against historical biomedical data and identifies novel biomedical insights |

The transformative potential of LLMs is further exemplified by their applications in addressing long-standing scientific challenges. For instance, AlphaFold Jumper et al. [2021] has revolutionized protein structure prediction, resolving key bottlenecks in drug discovery and expediting therapeutic innovation. Similarly, Crispr-GPT Huang et al. [2024] streamlines the design of gene-editing experiments, reducing the cognitive and procedural burdens on researchers while accelerating the pace of scientific advancement. In addition, MOLIERE Sybrandt et al. [2017, 2018]demonstrates how text mining and biomedical knowledge graphs can aid hypothesis validation by retrospectively testing hypotheses against historical data. These tools underscore the dual role of LLMs in augmenting human capabilities and enabling breakthroughs that transcend traditional boundaries of scientific inquiry. Table 1 provides a comprehensive overview of notable tools and their contributions to hypothesis generation and validation, illustrating the breadth of impact that these systems have achieved.

Despite their remarkable capabilities, LLMs exhibit fundamental differences from human researchers, particularly in reasoning and knowledge synthesis. Human cognition is characterized by inherent

intuition, creativity, and contextual understanding Langley and Simon [2013], enabling the pursuit of unconventional pathways that often lead to groundbreaking discoveries Fawzi et al. [2022]. Conversely, LLMs operate within the probabilistic framework of pattern recognition, constrained by the biases and limitations of their training datasets Tari et al. [2010]. While this allows for efficient data processing and knowledge reinforcement, it often results in the perpetuation of established ideas rather than the generation of genuinely novel concepts. This divergence underscores the necessity of harmonizing human intuition with machine-driven capabilities to achieve meaningful progress in scientific discovery Fok and Weld [2024], Zhou et al. [2024], Wang et al. [2023]. The success of agentic AI frameworks like SciAgents Ghafarollahi and Buehler [2024b] further demonstrates how hybrid human-AI collaboration can refine hypotheses, ensuring that machine-generated insights remain aligned with human reasoning. Such integration emphasizes the complementary roles of humans and machines, with each compensating for the limitations of the other.

Agentic AI systems Paul et al. [2024], White [2024], Chan et al. [2023], Sulc et al. [2024], Qiu and Lan [2024], powered by Large Language Models (LLMs), are reshaping the landscape of scientific discovery by automating routine but essential tasks such as data analysis, hypothesis formulation, and literature synthesis. These systems allow researchers to redirect their cognitive resources toward more creative and complex endeavors, thereby augmenting human ingenuity rather than replacing it. A recent comprehensive survey Gridach et al. [2025] provides a detailed taxonomy of agentic systems, distinguishing between autonomous and collaborative frameworks and categorizing their deployment across domains like chemistry, biology, and materials science. The study highlights their role across the full research lifecycle—from ideation and literature review to experimentation and scientific writing—and emphasizes human-AI collaboration and system calibration as pivotal directions for future development. Notable tools such as SciAgents Ghafarollahi and Buehler [2024b], AI Co-Scientist Gottweis et al. [2025], and reinforcement learning-driven frameworks for materials discovery Gruver et al. [2024] exemplify how domain-specific repositories and real-time feedback loops can enhance the contextual relevance and applicability of generated hypotheses. Likewise, systems like Discovering Faster Matrix Multiplication Algorithms Fawzi et al. [2022] and causal inference models in biomedical research Jha et al. [2019] further demonstrate the versatility of LLM-integrated scientific agents. However, despite these advancements, challenges remain, particularly concerning over-reliance on pre-existing data, limited novelty in generated hypotheses, and ethical considerations in high-stakes domains like healthcare Tang et al. [2024], Shavit et al. [2023]. These limitations underscore the need for responsible, human-in-the-loop designs that balance automation with domain expertise and ethical oversight.

This survey distinguishes itself by offering a holistic, interdisciplinary perspective on hypothesis creation and validation using LLMs and related AI systems. Unlike previous works that focus narrowly on specific domains such as biomedicine Sybrandt et al. [2017], Jumper et al. [2021] or materials science Jain et al. [2020], Gruver et al. [2024], this survey highlights the versatility of LLMs across a diverse array of fields, including social sciences Touvron et al. [2023], environmental studies Wang et al. [2023], and computational sciences Fok and Weld [2024]. By synthesizing insights from state-of-the-art studies and frameworks, this survey bridges the gap between theoretical advancements and practical applications Zhang et al. [2024a], Lu et al. [2024]. This comprehensive approach not only underscores the transformative potential of LLMs but also illuminates their adaptability to the multifaceted challenges of contemporary scientific research.

In addition to delineating the opportunities presented by LLMs, this work systematically identifies critical barriers to their effective utilization. Key challenges include enhancing the novelty of generated hypotheses, improving the feasibility of proposed ideas, mitigating data biases, and addressing the interdisciplinary adaptability of AI-driven methodologies. To overcome these obstacles, this survey proposes actionable strategies such as the development of generative exploration models Chen et al. [2024], the incorporation of human-in-the-loop systems Tari et al. [2010], and the fine-tuning of models for domain-specific applications. Furthermore, by emphasizing principles of transparency, fairness, and inclusivity, this survey addresses the ethical and practical considerations associated with deploying LLMs in scientific discovery Bommasani et al. [2021], Shavit et al. [2023], Fok and Weld [2024]. Through these contributions, this work provides a roadmap for harnessing the full potential of LLMs in advancing the frontiers of knowledge.

The remainder of this section outlines the prevailing approaches to scientific knowledge creation and validation, critiques their inherent limitations, and elaborates on the key contributions of this survey

paper. In doing so, it establishes a foundation for understanding the transformative impact of LLMs and related AI systems on the landscape of scientific discovery.

## 1.1  Overview of Current Approaches to Scientific Hypothesis Generation and Validation

The integration of LLMs into scientific research has initiated a paradigm shift in how hypotheses are generated and validated. Leveraging their capacity to process and synthesize vast volumes of domain-specific data, LLMs empower researchers to uncover latent patterns, generate insights, and explore relationships that are often inaccessible through traditional methodologies. Their scalability and computational efficiency make LLMs well-suited to address complex, interdisciplinary challenges across diverse scientific domains Zhang et al. [2024a], Lu et al. [2024]. LLMs have significantly transformed the hypothesis generation process by enabling data-driven exploration across structured and unstructured sources. Knowledge graph-based systems, such as MOLIERE Sybrandt et al. [2018, 2017] and SciAgents Ghafarollahi and Buehler [2024b], facilitate the discovery of novel connections by mapping semantic relationships in fields like biomedicine and materials science. Complementing these, retrieval-augmented generation (RAG) frameworks, exemplified by VELMA Schumann et al. [2024] and Chemist-X Chen et al. [2024], integrate curated knowledge bases with generative modeling to produce hypotheses that are both contextually grounded and creatively extended.

These methodologies are complemented by multi-omics integration platforms, such as VirtualPlant Katari et al. [2010] and BioLunar Wysocki et al. [2024], which synthesize genomic and pharmacological data to foster cross-disciplinary discoveries. In materials science, AI-driven hypothesis generation has been instrumental in predicting new material properties Gruver et al. [2024]. Advancing hypothesis generation, machine learning models equipped with reinforcement learning capabilities adapt dynamically to data-rich, iterative environments, as is commonly observed in drug discovery Blanco-Gonzalez et al. [2023], Tari et al. [2010]. Text mining tools, like Dyport Tyagin and Safro [2024], track the evolution of concepts across large textual datasets, enabling hypothesis creation in genomics and biomedical research. Other approaches, such as knowledge graph embeddings Wang et al. [2023], improve the scalability of hypothesis generation by structuring scientific knowledge into machine-readable formats. Together, these diverse methodologies illustrate how LLMs and AI-driven systems are reshaping the scientific discovery process, providing new pathways for interdisciplinary research and novel hypothesis generation.

Hypothesis validation is equally important to the scientific process, with LLMs playing a critical role in ensuring that proposed ideas are both scientifically plausible and testable. Tools like Labbench Laurent et al. [2024] and AgentClinic Schmidgall et al. [2024] facilitate experimental validation by simulating laboratory conditions to assess reproducibility and statistical significance. Simulation-based approaches provide cost-effective virtual environments for testing hypotheses in fields such as biomedicine Qi et al. [2024], Ghafarollahi and Buehler [2024a] and robotics Schumann et al. [2024]. Predictive model-based validation evaluates hypotheses using advanced metrics such as Bayesian posterior probability and prediction accuracy Jha et al. [2019], Chen et al. [2024]. These methodologies are augmented by cross-domain validation systems, which test the generalizability of hypotheses across diverse scientific domains Sybrandt et al. [2017], Zhou et al. [2024]. Additionally, crowdsourced validation platforms harness collective intelligence to evaluate hypotheses in social sciences and education, utilizing metrics like consensus scores Kim and Segev [2018], Touvron et al. [2023].

Multi-agent validation systems Ma et al. [2024] capitalize on distributed expertise to enhance collaborative validation processes, while causal inference frameworks Jha et al. [2019], Fok and Weld [2024] ensure the robustness of causal relationships through advanced structural causal modeling techniques. Dynamic validation tools Wang et al. [2023], Tang et al. [2024] incorporate real-time data to continuously refine validation processes, offering metrics like anomaly detection precision. Benchmarking platforms Qi et al. [2024], Aubin Le Quéré et al. [2024], Blanco-Gonzalez et al. [2023] ensure the reproducibility of validation efforts by measuring performance against standardized benchmarks. Moreover, iterative human-AI collaboration systems Shavit et al. [2023], Tari et al. [2010] combine human insights with AI-driven validation methodologies, enhancing explainability and user satisfaction.

The tools and datasets supporting these advancements form the foundation of contemporary knowledge creation and validation processes. Biomedical repositories like PubMed National Center for Biotechnology Information [2025], Gene Ontology (GO) Ashburner et al. [2000], and UK

Biobank Conroy et al. [2023] provide structured, high-quality data essential for hypothesis generation in areas such as genomics and drug discovery. Tools like MOLIERE Sybrandt et al. [2017] and Chemist-X Chen et al. [2024] effectively utilize these datasets to uncover novel connections by leveraging knowledge graphs and retrieval-augmented generation. Similarly, material science repositories, including the Materials Project Jain et al. [2020] and ChemBench Walker et al. [2010], offer comprehensive data on chemical compositions and properties, empowering researchers to hypothesize about new materials. AI-driven materials discovery frameworks, such as reinforcement learning-based methods Gruver et al. [2024], further enhance the predictive accuracy of material property estimations. Interdisciplinary repositories, hosted by platforms like Hugging Face[2], integrate multi-modal data to enable cross-domain hypothesis generation.

Simulation tools facilitate virtual testing environments, providing scalable validation methods for drug interactions, materials discovery, and robotics Qi et al. [2024], Schumann et al. [2024]. Meanwhile, open-source platforms like Hugging Face provide APIs and pre-trained models to streamline the integration of LLM capabilities into research workflows Touvron et al. [2023]. These platforms are increasingly adopted in both academic and industrial research, demonstrating their versatility in supporting large-scale scientific discovery. Through the synergistic application of these tools and methodologies, LLMs continue to advance scientific exploration by enabling more efficient and innovative approaches to hypothesis generation and validation. By grounding these processes in robust datasets and cutting-edge computational frameworks, researchers are well-positioned to uncover transformative insights that address some of the most pressing challenges of our time.

## 1.2   Challenges in Scientific Hypothesis Generation and Validation

Hypothesis creation, while significantly advanced by the capabilities of LLMs, faces several inherent limitations that impede the development of novel and impactful ideas. One critical challenge lies in the training paradigms of LLMs, which often replicate established knowledge patterns, thereby limiting the generation of truly innovative hypotheses. Techniques such as counterfactual reasoning and anomaly detection Fok and Weld [2024], Wang et al. [2023], Weidinger et al. [2021] are proposed to tackle the challenge. These techniques can be integrated into training processes to encourage deviation from conventional norms. Promoting novelty further requires the application of contrastive learning Hu et al. [2024] and dynamic retraining Han and Fu [1994] to explore uncharted scientific territories, especially in interdisciplinary contexts where cross-domain interactions hold the potential for groundbreaking insights Zhou et al. [2024], Chen et al. [2024]. However, achieving these insights requires using curated datasets and enhanced semantic mapping techniques to bridge disparate fields effectively Shavit et al. [2023], Jha et al. [2019]. The integration of structured knowledge graphs has been proposed to improve scientific hypothesis mapping by identifying non-obvious connections across domains Sybrandt et al. [2017]. Additionally, the scalability of data integration remains a persistent issue, with the need for scalable architectures and real-time synthesis to manage large and diverse datasets Touvron et al. [2023], Tang et al. [2024]. AI-driven hypothesis generation frameworks must also address limitations related to model adaptability, ensuring that evolving knowledge bases are incorporated into the reasoning process Gruver et al. [2024]. Another pressing concern is the interpretability of hypothesis generation systems, which often function as opaque "black boxes". Addressing this issue through logic-based reasoning and explainable AI (XAI) methodologies can significantly enhance trust and usability among researchers Shavit et al. [2023], Tari et al. [2010]. By integrating transparency-enhancing techniques such as causal inference models Jha et al. [2019] and retrieval-augmented reasoning Fawzi et al. [2022], LLMs can improve their ability to justify generated hypotheses in a scientifically rigorous manner.

Equally critical in hypothesis development, the validation phase ensures that proposed ideas are scientifically plausible and relevant. However, LLMs often struggle with domain-specific adaptability, as nuanced validation criteria vary significantly across specialized fields. Modular architectures incorporating domain-specific constraints offer a promising solution for improving validation accuracy Jha et al. [2019], Tang et al. [2024]. Another challenge is interdisciplinary validation, which requires adaptive systems capable of activating discipline-specific submodules to ensure the relevance of hypotheses across various fields Zhou et al. [2024], Sybrandt et al. [2017]. Furthermore, the reliance on computational metrics for validation often disconnects hypotheses from real-world feasibility. Simulation tools provide scalable virtual environments that can pre-test hypotheses effectively Qi

---

[2]https://huggingface.co/

et al. [2024], Schumann et al. [2024]. These methods have proven especially useful in biomedical research and materials science, where experimental verification can be resource-intensive Gruver et al. [2024]. However, high-risk, potentially transformative hypotheses are often penalized by conventional validation metrics that favor incremental advancements. To address this, risk-weighted evaluation frameworks must balance potential impact with empirical reliability Shavit et al. [2023], Fok and Weld [2024]. Additionally, multi-criteria validation approaches, combining metrics such as relevance, novelty, and feasibility, can provide a holistic assessment of hypothesis quality Qi et al. [2024], Aubin Le Quéré et al. [2024], Aubin Le Quéré et al. [2024]. Incorporating human-in-the-loop validation frameworks can further enhance hypothesis evaluation by integrating expert feedback into AI-driven validation systems Tari et al. [2010], Wang et al. [2023]. These methodologies contribute to a more rigorous and adaptable validation process, ensuring that LLM-generated hypotheses align with real-world applicability and scientific standards.

**Limitations in Achieving Novelty and Feasibility with LLMs.** The pursuit of novelty and feasibility in hypothesis creation is central to scientific progress, yet LLMs face fundamental challenges in both areas. A key issue lies in their training paradigm: LLMs are primarily trained on large corpora of existing, historically grounded data, which inherently biases them toward established knowledge patterns rather than fostering the generation of disruptive or paradigm-shifting insights Zhou et al. [2024], Ghafarollahi and Buehler [2024b]. This reliance results in a regression to the mean in idea generation, where LLMs tend to prioritize statistically likely continuations over epistemic risk-taking, even when prompted to be creative. Consequently, LLMs often generate variants of well-known ideas and rarely propose counterfactuals or unconventional hypotheses Fok and Weld [2024], Jha et al. [2019].

To overcome this limitation, several techniques have been proposed. Contrastive learning and generative exploration models have been developed to encourage semantic divergence and novelty Wang et al. [2023], Fawzi et al. [2022], Tang et al. [2024]. Reinforcement learning with novelty-seeking reward signals has also shown promise in promoting the exploration of low probability, high impact hypotheses Gruver et al. [2024], Blanco-Gonzalez et al. [2023]. However, these approaches remain experimental and require careful design to avoid compromising scientific rigor. In addition, novelty thresholds vary across disciplines, necessitating dynamic adjustment mechanisms within LLMs to align with domain-specific standards Zhou et al. [2024], Shavit et al. [2023]. Generating cross-disciplinary insights poses further complexity. Effective novelty requires semantic mapping across curated, multi-disciplinary datasets to uncover novel associations Ghafarollahi and Buehler [2024b], Sybrandt et al. [2017]. Yet, data bias and conservatism remain persistent obstacles. Models trained on historical data often struggle to produce original ideas, reinforcing conventional thinking. Diversified training corpora and novelty-boosting algorithms have been proposed to mitigate these effects Chen et al. [2024], Wang et al. [2023]. Additionally, surface-level similarity metrics are insufficient for detecting deeper conceptual innovations, underscoring the need for models capable of identifying unique theoretical implications Gruver et al. [2024]. Retrieval augmented reasoning can enhance this capability by grounding hypotheses in diverse, relevant contexts Jha et al. [2019].

Feasibility presents a parallel set of challenges. LLMs often generate hypotheses without considering practical constraints, limiting their real-world applicability. Effective feasibility requires grounding in empirical evidence and domain-specific constraints. Techniques such as the integration of multi-modal data, including experimental results and sensor outputs, have been proposed to improve feasibility assessment Qi et al. [2024], Ghafarollahi and Buehler [2024b]. Foundation models with physical interaction capabilities can further bridge the gap between theory and experimental validation Blanco-Gonzalez et al. [2023], Jha et al. [2019]. For example, models integrated with robotic platforms, such as those used in automated laboratories or robotic chemists King [2011], Fakhruldeen et al. [2022], can physically conduct experiments based on model-generated hypotheses, enabling real-time feedback and refinement. These systems, like the "robot scientist" Eve Williams et al. [2015] used in drug discovery, embody physical interaction by selecting compounds, operating lab equipment, and analyzing results, closing the loop between hypothesis generation and empirical testing. In data-scarce scientific domains, synthetic data generation and few-shot learning can supplement existing datasets to improve generalization Wang et al. [2023], Chen et al. [2024]. Embedding field-specific constraints into LLM pipelines, such as resource availability or laboratory feasibility, can ensure that generated hypotheses are actionable Gruver et al. [2024], Shavit et al. [2023]. Simulation environments like ChemBench provide scalable virtual platforms for early-stage hypothesis testing Walker et al. [2010]. In biomedical and materials science, digital twin simulations offer higher fidelity

feasibility assessments by emulating real-world experimental settings Fawzi et al. [2022], Tang et al. [2024]. Tailoring feasibility metrics to reflect domain-specific requirements remains essential for producing hypotheses that are both practical and impactful Aubin Le Quéré et al. [2024], Touvron et al. [2023].

**Ethical Concerns in LLM-Driven Hypothesis Generation and Validation.** As LLMs gain autonomy in scientific research, several ethical challenges emerge that impact the trustworthiness, accountability, and inclusivity of AI-assisted discovery. First, LLMs trained on biased data may reproduce and even amplify social, demographic, and epistemic inequities, potentially marginalizing underrepresented perspectives or reinforcing dominant paradigms Shavit et al. [2023], Weidinger et al. [2021]. This risk is particularly serious in fields like biomedicine and public policy, where disparities in training data can lead to flawed or discriminatory hypotheses. Second, the phenomenon of AI hallucination Venkit et al. [2024] presents a threat to scientific integrity. LLMs may generate outputs that appear fluent and scientifically plausible but are factually incorrect or unsupported by evidence. Without transparent reasoning or source attribution, these outputs can mislead researchers and corrupt the hypothesis evaluation pipeline Tang et al. [2024]. This is especially dangerous when outputs are integrated into high-stakes domains such as clinical decision-making or environmental policy. Third, accountability is difficult to establish when LLM-generated hypotheses lead to errors or unintended consequences. Whether responsibility should lie with model developers, research users, or deploying institutions remains unclear. In the absence of mechanisms for provenance tracking, version control, and responsible usage documentation, resolving these questions remains difficult Jaradeh et al. [2019]. Ethical safeguards must be embedded throughout the hypothesis generation and validation pipeline to mitigate these concerns. These include explainable AI (XAI) methods, allowing researchers to understand and verify model reasoning, and human-in-the-loop frameworks integrating expert oversight into model output evaluation Jaradeh et al. [2019], Tang et al. [2024]. Auditing protocols and model cards Kennedy-Mayo and Gord [2025], Nunes et al. [2024] can also help disclose ethical risks, intended use cases, and model limitations, promoting transparency and accountability.

**Regulatory and Policy Implications.** The increasing integration of LLMs into scientific workflows presents a growing governance challenge. Despite their widespread use, few regulatory standards exist for evaluating or auditing AI-generated hypotheses. Traditional scientific metrics do not adequately capture emerging concerns introduced by LLMs, such as ethical risk, reproducibility, and explainability. To ensure responsible deployment, domain-specific policy frameworks are needed to support traceable, verifiable, and equitable use of LLMs in science. These frameworks should include auditing tools to verify provenance, model documentation standards, and transparent disclosures of AI contributions in research publications. Open science principles must be preserved by promoting access to open-source models and datasets, especially for researchers in low-resource environments. Institutions, funding agencies, and regulatory bodies should adapt review criteria and funding guidelines to account for the role of LLMs in the research process. This includes evaluating whether proper oversight, validation mechanisms and ethical safeguards are in place. Without these changes, scientific discovery risks becoming dependent on opaque systems that operate without clear accountability or alignment with community norms.

## 1.3 Contributions of this Study

This survey paper comprehensively examines the challenges and opportunities in automating hypothesis creation and validation using LLMs. As LLMs increasingly contribute to scientific research across diverse disciplines, understanding their limitations and potential is critical for fostering innovation. Despite their transformative capabilities, LLMs face significant challenges in generating novel, feasible, and impactful hypotheses. By consolidating methodologies, identifying gaps, and proposing actionable strategies, this survey equips the research community with a roadmap to enhance the effectiveness of LLMs in hypothesis generation and validation.

The primary contributions of this study are as follows:

- **Comprehensive Review of Current Approaches:** This survey offers a structured, interdisciplinary overview of current approaches to hypothesis creation and validation. It synthesizes insights from diverse scientific fields, bridging gaps in the literature and providing a unified perspective to guide future research directions.
- **Identification of Challenges and Gaps:** Key challenges in hypothesis creation and validation, such as promoting novelty, improving feasibility assessments, and addressing data

biases, are systematically analyzed. These insights illuminate where LLMs fall short and provide clarity for researchers and practitioners aiming to develop robust and innovative systems.

- **Interdisciplinary Relevance:** Highlighting the versatility of LLMs, this survey demonstrates their applicability across fields such as biomedicine, materials science, social sciences, and environmental research. This survey illustrates how LLMs can adapt to varied scientific contexts and foster innovation across disciplines by showcasing use cases and domain-specific challenges.

- **Establishing a Framework for Ethical and Practical Applications:** Ethical and practical considerations are emphasized to ensure that LLMs are deployed responsibly and effectively. This survey sets foundational guidelines for creating transparent, fair, and inclusive systems, fostering trust and usability in scientific workflows.

- **Actionable Insights for Transformative Advancements:** By encouraging a shift from static, one-size-fits-all approaches to dynamic, adaptive, and interdisciplinary systems, this survey aligns with the evolving demands of scientific research. It highlights how data-driven methodologies can redefine the pace and scope of discovery, providing tools and insights for transformative advancements in hypothesis creation and validation.

- **Proposal of Novel Strategies and Future Directions:** This survey presents forward-looking strategies to overcome identified challenges. These include generative exploration models, hybrid human-AI systems, and risk-tolerant validation frameworks that equip researchers and developers with practical tools to enhance LLM-driven scientific exploration.

In summary, this survey addresses the critical need for a holistic, interdisciplinary resource that integrates theoretical advancements with practical applications. By identifying gaps, offering actionable strategies, and emphasizing ethical considerations, to empower researchers, practitioners, and policymakers to harness the transformative potential of LLMs for scientific innovation. This survey offers a foundation for exploring approaches to hypothesis creation and validation, contributing to the development of data-driven scientific discovery.

## 2 Definitions and Overview

This section introduces key concepts essential to this survey and outlines the structure of the paper, providing a roadmap for the topics discussed in subsequent sections.

### 2.1 Definitions

**Hypothesis:** A hypothesis is a tentative explanation, relationship, or proposition that can be empirically tested or theoretically evaluated. In scientific research, hypotheses serve as foundational units for exploration, guiding the formulation of experiments or analytical tasks. Within the context of LLM-driven scientific discovery, a hypothesis can range from a declarative statement proposing a causal relationship (e.g., "Gene X inhibits Protein Y") to an abstract proposition extracted or synthesized from unstructured data (e.g., text-based summaries or concept clusters).

**Novelty:** Novelty in hypothesis creation Witt [2009], Hallsworth et al. [2023] refers to the degree to which a generated hypothesis differs from existing knowledge. It indicates the originality of the hypothesis by measuring its divergence from known concepts, relationships, or patterns in a given domain. Quantifying novelty is essential for evaluating whether hypotheses contribute meaningful advancements rather than reiterating existing ideas Zhou et al. [2024], Fok and Weld [2024], Shibayama et al. [2021].

Novelty can be quantified using similarity or distance metrics. Given a generated hypothesis $H$ represented as a vector $h$ in an embedding space, and a set of existing hypotheses $H_i$ with vectors $h_i$, the novelty $N(H)$ of $H$ can be computed as the inverse of the average cosine similarity between $h$ and each $h_i$:

$$N(H) = 1 - \frac{1}{|S|} \sum_{i \in S} \text{cosine\_similarity}(h, h_i), \tag{1}$$

where $S$ is the set of existing hypotheses.

High novelty implies *low similarity* to existing knowledge, suggesting *originality* and a potential breakthrough or unexplored idea. However, ensuring novelty while maintaining scientific validity is a challenging balance, as excessive deviation from known knowledge can lead to implausible or untestable hypotheses Gruver et al. [2024], Shavit et al. [2023]. Contrastive learning and retrieval-augmented reasoning have been proposed as strategies to improve novelty in LLM-generated hypotheses by refining semantic divergence while preserving logical consistency Kim et al. [2024].

**Feasibility:** Feasibility Walker [1987], Song et al. [2024] assesses whether a generated hypothesis is practically testable and grounded within the known constraints of the domain. It evaluates the likelihood that a hypothesis can be experimentally validated or implemented in real-world settings. Feasibility is particularly critical in domains such as biomedicine and materials science, where empirical validation requires substantial experimental resources Wang et al. [2024a]. The feasibility $F(H)$ of a hypothesis $H$ can be calculated as a weighted combination of empirical and theoretical scores. If $f_{empirical}$ represents empirical feasibility, for instance, the availability of necessary data, equipment, or methods, and $f_{theoretical}$ represents the theoretical validity based on domain knowledge, the feasibility is defined as:

$$F(H) = w_{emp} \cdot f_{empirical} + w_{theo} \cdot f_{theoretical}, \tag{2}$$

where $w_{emp}$ and $w_{theo}$ are weights reflecting the importance of each factor, subject to $w_{emp} + w_{theo} = 1$. A high feasibility score indicates that the hypothesis is both theoretically sound and practically achievable, increasing its potential value for real-world testing Shavit et al. [2023], Jha et al. [2019].

Multi-modal feasibility assessments, incorporating experimental data, sensor inputs, and real-world simulations, have been proposed to improve LLM-driven hypothesis validation Fawzi et al. [2022], Chen et al. [2024]. AI-powered frameworks such as digital twin simulations Tang et al. [2024] and human-in-the-loop validation further enhance feasibility by dynamically refining hypotheses based on real-time feedback. Tailoring feasibility metrics to discipline-specific challenges ensures that LLM-generated hypotheses remain actionable and impactful in scientific research Aubin Le Quéré et al. [2024], Touvron et al. [2023].

**Open Domain:** An open domain in hypothesis generation refers to a broad, unrestricted field where generated hypotheses are not confined to specific topics or predefined structures. In open-domain settings, hypotheses may span across multiple fields, incorporating interdisciplinary knowledge Zhou et al. [2024], Ghafarollahi and Buehler [2024b]. The open domain can be represented as a large, unrestricted hypothesis space $\mathcal{H}_{open}$, where any hypothesis $H$ is possible. It is formally defined as:

$$\mathcal{H}_{open} = \{H \mid H \in \text{Any Topic or Field}\}. \tag{3}$$

Open-domain hypothesis generation requires models that can generalize across diverse fields, enabling novel, interdisciplinary insights, but often making validation and relevance assessment more challenging Sybrandt et al. [2018], Wang et al. [2023]. Techniques such as retrieval-augmented generation (RAG) and knowledge graph-based reasoning have been employed to structure open-domain hypothesis generation, ensuring hypotheses remain scientifically plausible while preserving creativity Chen et al. [2024], Jha et al. [2019]. One of the key challenges in open-domain hypothesis generation is contextual grounding, as LLMs may generate syntactically correct hypotheses but lack empirical feasibility Shavit et al. [2023], Aubin Le Quéré et al. [2024]. To address this, adaptive fine-tuning methods that incorporate domain constraints dynamically have been proposed to improve the relevance of generated hypotheses while maintaining the breadth of interdisciplinary exploration Tang et al. [2024], Touvron et al. [2023].

**Closed Domain:** A closed domain refers to a specific, restricted area of knowledge in which generated hypotheses are limited to a particular subject or closely related subfields. In a closed-domain setting, hypotheses are generated with narrowly defined constraints, making them more easily verifiable within the domain Qi et al. [2024], Gruver et al. [2024]. Closed domain can be represented by a bounded hypothesis space $\mathcal{H}_{closed}$, restricted by a set of domain constraints $C$:

$$\mathcal{H}_{closed} = \{H \mid H \text{ adheres to constraints } C\}, \tag{4}$$

where $C$ includes specific requirements or parameters related to the domain (e.g., biomedicine, materials science). Closed-domain hypothesis generation enhances relevance and ease of validation but may limit the scope for cross-disciplinary or radically innovative insights Jha et al. [2019], Shavit et al. [2023]. To improve domain-specific accuracy, retrieval-augmented generation (RAG) and fine-tuned foundation models have been employed in closed-domain settings Chen et al. [2024], Wang

et al. [2023]. Tools such as MOLIERE for biomedical hypothesis validation Sybrandt et al. [2018] and SciAgents for structured scientific discovery Ghafarollahi and Buehler [2024b] demonstrate the efficacy of constrained knowledge generation. However, over-restriction in closed-domain settings may reduce the exploratory potential of LLMs, necessitating hybrid approaches that allow limited cross-domain adaptation while maintaining domain constraints Tang et al. [2024], Touvron et al. [2023].

**Relevance of Generated Hypothesis:** Relevance assesses how well a hypothesis aligns with the target scientific community's goals, needs, or pressing questions. It measures contextual importance within the domain Wang et al. [2023]. Relevance can be computed via citation-based retrieval scores, topic overlap, or alignment with funding priorities.

**Quality of Generated Hypothesis:** The quality of a generated hypothesis reflects its potential scientific impact, evaluated through a combination of novelty, feasibility, and relevance. A high-quality hypothesis introduces new achievable and significant insights within the scientific domain. The quality $Q(H)$ of a hypothesis $H$ can be expressed as a weighted aggregate of novelty $N(H)$, feasibility $F(H)$, and relevance $R(H)$:

$$Q(H) = w_N \cdot N(H) + w_F \cdot F(H) + w_R \cdot R(H), \qquad (5)$$

where $w_N$, $w_F$, and $w_R$ are weights assigned based on the importance of each component for the specific research objective, and $w_N + w_F + w_R = 1$. A hypothesis with high-quality scores well across novelty, feasibility, and relevance, making it a strong candidate for empirical testing and potential scientific advancement Wang et al. [2023], Jha et al. [2019]. Multi-criteria evaluation frameworks have been proposed to improve hypothesis quality assessment, integrating both quantitative and qualitative metrics Shavit et al. [2023], Chen et al. [2024]. For instance, contrastive learning techniques enable LLMs to refine hypotheses by balancing novelty and plausibility, reducing the risk of generating implausible or irrelevant ideas Hu et al. [2024]. Similarly, human-in-the-loop validation allows domain experts to iteratively refine LLM-generated hypotheses, ensuring that AI-driven scientific discoveries align with empirical research priorities Aubin Le Quéré et al. [2024], Touvron et al. [2023].

**Bad Hypothesis:** A bad hypothesis is one that is either trivial, incorrect, non-novel, or infeasible. Such a hypothesis lacks scientific merit due to redundancy with existing knowledge, conceptual errors, or impracticality for empirical testing. The criteria for categorizing a generated hypothesis as bad include:

1. *Low Novelty:* Highly similar to known hypotheses, offering little or no new information Qi et al. [2024].

2. *Low Feasibility:* Practically or theoretically untestable, making empirical validation unrealistic Kim and Segev [2018].

3. *Conceptually Incorrect:* These hypotheses violate domain logic or introduce contradictions. They are implicitly penalized via low $f_{theoretical}$ in feasility and can also be flagged using symbolic consistency checks or NLI-based verification frameworks Zhou et al. [2024], Sybrandt et al. [2018].

Let $N(H)$ denote novelty, $F(H)$ feasibility, and $E(H)$ represent correctness (with $E(H) = 1$ for a valid hypothesis and 0 for a conceptually flawed hypothesis). A hypothesis $H$ is classified as bad if:

$$N(H) < \tau_N \quad \text{or} \quad F(H) < \tau_F \quad \text{or} \quad E(H) = 0, \qquad (6)$$

where $\tau_N$ and $\tau_F$ are thresholds for novelty and feasibility, respectively. These thresholds are domain-specific. For instance, earlier studies Sybrandt et al. [2018], Tang et al. [2024] employ percentile-based thresholds derived from baseline distributions, whereas others adopt tunable cutoffs informed by downstream validation outcomes or expert evaluations. Identifying and filtering bad hypotheses can prevent resource wastage in experimental stages, allowing focus on hypotheses with potential scientific value Beltagy et al. [2019].

LLMs can be enhanced through knowledge-driven hypothesis generation that leverages structured data sources and domain expertise to mitigate the generation of bad hypotheses. Additionally, retrieval-augmented validation frameworks such as SciFact and MOLIERE improve hypothesis filtering by cross-referencing scientific literature. AI-driven hypothesis verification techniques, integrating NLP-based consistency checks, further enhance the rejection of incorrect or redundant hypotheses.
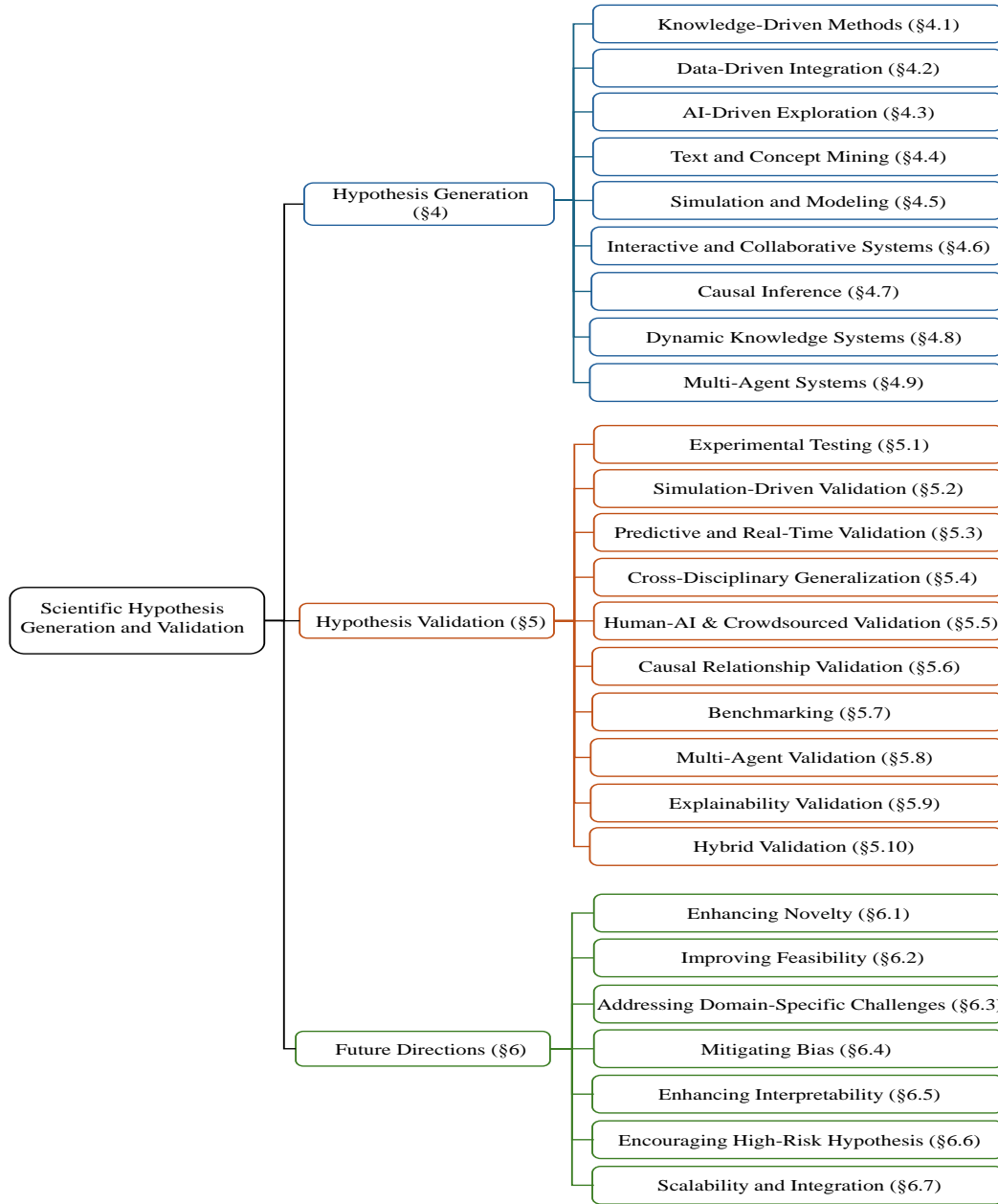
Figure 2: Flow diagram of the survey structure. This figure guides the reader through the organization of the paper, beginning with hypothesis creation approaches (§4), progressing through hypothesis validation methods (§5), and culminating in open challenges and future directions (§6). It highlights how various components of scientific discovery, from AI-driven exploration to validation frameworks, are interconnected in the survey's narrative.

## 2.2 Overview and Structure of the Paper

This survey provides a systematic and comprehensive exploration of scientific knowledge creation and validation using Large Language Models (LLMs). It is organized into sections that address the methodologies, challenges, and opportunities in advancing this field. The structure is designed to guide readers through an understanding of the state-of-the-art, the limitations of current systems, and the opportunities for future innovation. Figure 2 provides an overview of the survey structure.

- **Datasets Supporting Hypothesis Generation and Validation (Section 3):** Data serves as the foundation for hypothesis generation and validation. Section 3 describes available datasets (e.g., SciFact, PubMedQA, AVeriTeC) and tools like MOLIERE, ChemBench, and Project Jupyter that support hypothesis analysis and validation. It highlights the importance of understanding the role of datasets, their domain-specific applications, and the challenges of ensuring data quality, fairness, and accessibility to enable meaningful scientific exploration.

- **Hypothesis Generation Approaches (Section 4):** Section 4 explores the diverse methodologies employed for hypothesis generation, emphasizing the need to understand how LLMs identify patterns, uncover relationships, and generate new ideas. It highlights the importance of categorizing approaches to assess their strengths and limitations in addressing various scientific challenges.

- **Hypothesis Validation Approaches (Section 5):** Validation is a critical step in ensuring that hypotheses are scientifically plausible and impactful. Section 5 discusses the necessity of robust validation methods to evaluate hypotheses' feasibility, novelty, and relevance while addressing the constraints posed by different scientific contexts.

- **Future Directions for Advancing the Field (Section 6):** Building on the identified challenges, Section 6 motivates the need for actionable strategies to enhance the capabilities of LLMs. It discusses the potential of integrating advanced methodologies and fostering interdisciplinary collaboration to push the boundaries of hypothesis creation and validation.

- **Conclusion (Section 7):** The paper concludes in Section 7 by reflecting on the transformative potential of LLMs in scientific discovery. It emphasizes the urgency of adopting innovative approaches to maximize their impact and advocates for a shift toward systems that prioritize creativity, adaptability, and interdisciplinary applicability.

## 3 Datasets

Datasets are the cornerstone of hypothesis creation and validation, providing the foundational information required to generate, test, and refine hypotheses. Their selection influences not only the novelty and feasibility of hypotheses but also their scope and applicability across domains. From biomedicine and materials science to social science and artificial intelligence, datasets reflect the diverse nature of scientific inquiry and drive advancements by enabling tailored hypothesis generation and rigorous validation.

This section reviews key datasets widely used across scientific fields, categorizing them by structure, domain, and functionality. By exploring curated knowledge graphs, textual corpora, and specialized scientific data repositories, we highlight their role in supporting innovation, addressing complex challenges, and fostering interdisciplinary research. The systematic analysis of these datasets provides insights into their characteristics, evaluation metrics, and potential for enabling groundbreaking discoveries.

**PubMed National Center for Biotechnology Information [2025]:** PubMed is a foundational dataset in biomedical research, featuring over 34 million abstracts and citations across clinical medicine, pharmacology, and molecular biology. It supports hypothesis creation and validation by enabling co-occurrence analysis, uncovering relationships between biomedical entities, and facilitating natural language processing tasks such as entity recognition and document classification. Key evaluation metrics include co-occurrence analysis for identifying term relationships, ROC-AUC for assessing model performance in pattern recognition, and relevance metrics for aligning abstracts with hypotheses. These tools ensure robust hypothesis testing and data-driven insights. PubMed fosters novelty by revealing previously unlinked relationships while providing evidence-based validation through its peer-reviewed content, making it indispensable for biomedical research.

Table 2: Summary of Datasets for Hypothesis Creation and Validation (Part 1)

| Dataset Name | Description | Statistics | Domain | Evaluation Metrics | Modality | Novelty (Y/N) | Feasibility (Y/N) |
|---|---|---|---|---|---|---|---|
| PubMed Abstracts National Center for Biotechnology Information [2025] | Biomedical literature database for hypothesis generation in biology and medicine | Over 34 million abstracts | Biomedicine | Relevance, Co-occurrence Analysis, ROC-AUC | Text | Y | Y |
| MeSH U.S. National Library of Medicine [2025] | Medical Subject Headings for categorizing PubMed content | 27,883 descriptors | Biomedicine | Similarity Metrics, Novelty Scoring | Text, Structured Metadata | Y | Y |
| ChEMBL Zdrazil et al. [2023] | Bioactive molecule database for drug discovery | Over 2 million compounds | Chemistry | Molecular Similarity, Drug-Likeness Scores | Text, Numerical | Y | Y |
| GENIA Corpus Kim et al. [2003] | Annotated biomedical text corpus for NLP tasks | Over 2,000 abstracts annotated with biomedical terms | Biomedicine | Precision, Recall, F1-Score | Text | N | N |
| Open Graph Benchmark (OGB) Hu et al. [2020] | Large-scale graph datasets for hypothesis testing on network data | Over 100 datasets for benchmarking graph-based tasks | AI, Graph Analysis | Graph Accuracy, Link Prediction Accuracy | Graphs | Y | Y |
| UK Biobank Conroy et al. [2023] | Longitudinal dataset linking genetic and phenotypic data | Data from over 500,000 participants | Genomics | Correlation, Causal Inference | Text, Numerical, Genomic | N | Y |
| MATBench Dunn et al. [2020] | Material property prediction dataset for materials discovery | Includes over 100,000 material samples | Materials Science | Prediction Accuracy, RMSE | Numerical, Structured Data | Y | Y |
| ClimateNet Prabhat et al. [2020] | Dataset for climate change research and hypothesis validation | Over 10 years of climate observation data | Environmental Science | Temporal Trends, Anomaly Detection | Text, Numerical, Satellite Imagery | N | N |
| COCO Dataset Lin et al. [2014] | Dataset for image captioning and computer vision tasks | Over 300,000 images with captions | AI, Vision-Language Integration | Precision, BLEU Score | Image, Text | N | N |
| Gene Ontology (GO) Ashburner et al. [2000] | Hierarchical vocabulary for annotating gene products, supporting hypothesis generation and validation in genomics and molecular biology | Over 44,000 terms across three domains: biological process, molecular function, and cellular component | Genomics, Molecular Biology | Semantic Similarity, Graph-Based Validation, Annotation Consistency | Text, Graphs | N | Y |
| AHTech Electrolyte Additive Dataset Lin et al. [2025] | High-throughput electrochemical screening data for electrolyte additives in aqueous zinc batteries | 180 candidates, 200-cycle efficiency per sample | Electro-chemistry, Energy Storage | Coulombic Efficiency, Additive Ranking, Feature Correlation | Tabular, Text | Y | Y |
| CSKG-600 Borrego et al. [2025] | Expert-labeled hypothesis triples from a scholarly knowledge graph | 600 hypothesis statements with expert validation | Knowledge Graphs, Scholarly AI | Precision@K, Link Prediction Accuracy, Expert Agreement | Graphs, Text | Y | Y |

**MeSH (Medical Subject Headings) U.S. National Library of Medicine [2025]:** MeSH is a structured controlled vocabulary that categorizes biomedical content in databases such as PubMed, providing a hierarchical framework for precise hypothesis generation and validation. It facilitates the categorization of biomedical content, linking terms to established descriptors to support hypothesis validation and enabling similarity and novelty scoring for hypothesis evaluation. Key metrics include semantic similarity, which quantifies relationships using ontological structures, novelty scoring to measure the uniqueness of hypotheses, and term coverage to assess alignment with MeSH categories. MeSH supports novel hypothesis creation by enabling semantic exploration and unique term connections, while its hierarchical structure ensures reliability in hypothesis validation.

**ChEMBL Zdrazil et al. [2023]:** ChEMBL is a comprehensive database of bioactive molecules designed for hypothesis generation and validation in drug discovery and medicinal chemistry. It

Table 3: Summary of Datasets for Hypothesis Creation and Validation (Part 2)

| Dataset Name | Description | Statistics | Domain | Evaluation Metrics | Modality | Novelty (Y/N) | Feasibility (Y/N) |
|---|---|---|---|---|---|---|---|
| DrugBank Wishart et al. [2018] | Comprehensive dataset of drug-target interactions for drug discovery | Over 14,000 drugs and 6,000 protein targets | Pharmacology | Interaction Accuracy, Drug Efficacy Scores | Text, Structured Tables | Y | Y |
| AI2 Science Questions Clark et al. [2018] | Dataset for hypothesis testing in question answering and reasoning | Over 10,000 multiple-choice science questions | Education, AI Reasoning | Accuracy, Explainability Metrics | Text, Structured QA Format | Y | N |
| Materials Project Jain et al. [2013] | Database of materials properties and structures | Over 133,500 materials | Materials Science | Structural Matching, Novelty Filtering | Numerical, Structured Data | Y | Y |
| KEGG Pathway Kanehisa and Goto [2000] | Database of metabolic and signaling pathways | 540 pathways across species | Biomedicine, Genomics | Pathway Enrichment, Graph Metrics | Text, Graphs | Y | Y |
| American Community Survey (ACS) Bureau [2025] | Annual survey capturing demographic and social data in the US | 2.5 million responses/year | Social Science | Statistical Analysis, Demographic Comparisons | Text, Numerical | N | Y |
| Patent Data (USPTO) Patent and Office [2025] | Text and metadata of granted patents | Over 10 million patents | Technology, Innovation | Citation Analysis, Novelty Score | Text, Structured Metadata | Y | N |
| XSum Narayan et al. [2018] | Summarization dataset with diverse topics for NLP research | 226,711 summaries | NLP/Text Mining | BLEU, ROUGE for validation | Text | Y | Y |
| Cosmic Bamford et al. [2004] | Somatic mutation data for cancer research | 30,000 genes, 2 million mutations | Cancer Genomics | Mutation Analysis, Pathway Mapping | Text, Numerical | N | Y |
| Open Research Knowledge Graph (ORKG) Jaradeh et al. [2019] | Structured knowledge graph of research contributions | 3 million triples | Multidisciplinary | Graph Centrality, Novelty Detection | Text, Graphs, Structured Metadata | Y | Y |

facilitates the exploration of drug-target interactions, validates predictive models for structure-activity relationships (SAR), and enables analysis of compound efficacy and bioactivity. The dataset includes over 2 million molecules with bioactivity data, covering more than 14,000 biological targets. Metrics such as molecular similarity, drug-likeness scores, and binding affinity prediction accuracy support novel drug candidates and validate bioactivity insights. ChEMBL's experimental data ensures both novelty and feasibility in hypothesis testing.

**GENIA Corpus Kim et al. [2003]:** The GENIA Corpus is an annotated biomedical text dataset tailored for NLP tasks, facilitating hypothesis generation and validation in text mining and entity recognition. Comprising over 2,000 abstracts annotated with more than 36,000 unique terms, it focuses on biomedical entities like proteins and genes, emphasizing transcription factors and cellular signaling. Its primary use is benchmarking NLP models, with evaluation metrics such as precision, recall, and F1-score for entity recognition. While it excels in validating NLP techniques, its primary role is not novel hypothesis generation.

**Open Graph Benchmark (OGB) Hu et al. [2020]:** The Open Graph Benchmark (OGB) is a collection of over 100 graph datasets across domains like biology, chemistry, and computer science, designed for benchmarking machine learning models in graph-based tasks. It supports hypothesis generation about structural patterns in networks and validates graph-based hypotheses using metrics like graph accuracy, link prediction accuracy, and clustering coefficients. OGB provides predefined training, validation, and testing splits, making it ideal for exploring novel structural patterns and systematically validating hypotheses.

**UK Biobank Conroy et al. [2023]:** The UK Biobank offers a vast dataset of genetic, phenotypic, and health-related information from 500,000 participants, supporting hypothesis generation in genomics and personalized medicine. Covering over 800 phenotypic traits and 96 million genetic variants,

it enables causal inference and phenotypic prediction, evaluated through metrics like correlation coefficients, causal inference metrics, and prediction accuracy. While it excels in hypothesis validation, its focus is primarily on validating existing hypotheses rather than generating novel ones.

**MATBench Dunn et al. [2020]:** MATBench is a benchmark dataset for material property prediction, featuring over 100,000 material samples across categories like alloys, ceramics, and polymers. It supports hypothesis generation about material properties and validates predictive models using metrics like prediction accuracy, RMSE, and $R^2$. Its well-defined training, validation, and testing splits make it a robust resource for exploring novel material properties and systematically validating machine learning models in materials science.

**ClimateNet Prabhat et al. [2020]:** ClimateNet is an expert-labeled dataset designed for climate change research and hypothesis validation. It facilitates the identification and analysis of extreme weather patterns, enabling hypothesis generation in climate science by providing structured observational data. The dataset supports hypothesis validation by offering labeled climate events that can be used to train and test predictive models for extreme weather forecasting and climate anomaly detection. Key evaluation metrics include temporal trends analysis to assess long-term climate variations, anomaly detection for identifying deviations from expected climate behaviors, and spatial correlation metrics to measure the consistency of climate phenomena across regions. While ClimateNet is primarily used for validation rather than novel hypothesis discovery, its expert-curated labels and structured data representation make it a valuable tool for enhancing climate prediction models and improving the understanding of atmospheric processes.

**COCO Dataset Lin et al. [2014]:** The COCO (Common Objects in Context) dataset is a large-scale dataset designed for image captioning, object detection, and vision-language integration tasks. It enables hypothesis generation in artificial intelligence by providing a diverse set of annotated images for evaluating computer vision and natural language processing models. For example, a researcher might hypothesize that "transformer-based vision-language models generate more contextually accurate image captions than RNN-based models when multiple objects co-occur in complex scenes." This hypothesis can be tested using COCO's richly annotated images and corresponding captions. Hypothesis validation is supported through well-defined image-to-text relationships, allowing researchers to assess model performance in object recognition, scene understanding, and multimodal learning. Key evaluation metrics include precision and recall for object detection accuracy, BLEU score for measuring the alignment between generated and reference captions, and segmentation accuracy for validating instance-level object identification. While COCO primarily serves as a benchmarking dataset, its rich annotations and large-scale diversity facilitate advancements in AI-driven image interpretation and vision-language research.

**Gene Ontology (GO) Ashburner et al. [2000]:** The Gene Ontology (GO) dataset offers a hierarchical vocabulary for annotating gene products, supporting hypothesis generation and validation in genomics and molecular biology. It enables hypothesis generation by detailing gene functions, processes, and cellular components, and validates computational models for gene function prediction and pathway analysis. Key metrics include semantic similarity to measure functional similarities, graph-based validation using metrics like connectivity and path length, and annotation consistency to evaluate reliability. While GO is primarily used for hypothesis validation rather than novel discovery, its detailed annotations and hierarchical structure make it a robust tool for validating existing hypotheses.

**AHTech Electrolyte Additive Dataset Lin et al. [2025]:** This dataset was introduced as part of the AHTech platform for accelerating electrochemical discovery. It contains high-throughput screening data from 180 small-molecule electrolyte additives tested for aqueous zinc metal batteries. Each additive was characterized across 200 electrochemical cycles to determine Coulombic efficiency, enabling the training of machine learning models for additive performance prediction. The dataset supports hypothesis validation by uncovering structure-performance relationships using techniques such as Shapley Additive Explanations (SHAP) and Spearman correlation. Its high experimental fidelity and structured annotations make it valuable for data-driven hypothesis generation in electrochemistry and battery materials research.

**CSKG-600 Borrego et al. [2025]:** CSKG-600 is a benchmark dataset introduced to evaluate hypothesis generation over scholarly knowledge graphs. It consists of 600 candidate hypotheses manually labeled by domain experts as valid or invalid, supporting developing and evaluating link prediction models for scientific discovery. The dataset integrates structured triples with semantic and bibliometric metadata, enabling robust benchmarking of systems like ResearchLink. It facilitates

validation through ranking-based metrics such as Precision@K and domain-specific agreement scores. As one of the first domain-independent resources in this space, CSKG-600 is well-suited for hypothesis validation tasks involving multi-modal, interdisciplinary scientific knowledge.

**DrugBank Wishart et al. [2018]:** DrugBank integrates comprehensive data on drugs and molecular targets, facilitating hypothesis generation and validation in pharmacology and bioinformatics. With over 14,000 drugs and 6,000 protein targets, it supports predictions in drug efficacy and adverse effects, evaluated through metrics like interaction accuracy, RMSE for binding affinity predictions, and drug-likeness scores. DrugBank excels in discovering novel drug candidates and provides experimental and clinical data for robust validation.

**AI2 Science Questions Clark et al. [2018]:** The AI2 Science Questions dataset comprises over 10,000 multiple-choice questions, testing AI reasoning and knowledge representation in educational research. Covering topics from physics to earth science, it supports hypothesis generation about reasoning capabilities and evaluates model performance using metrics like accuracy, explainability scores, and logical consistency metrics. While it supports exploration of novel reasoning strategies in AI, it primarily focuses on evaluating reasoning rather than hypothesis feasibility.

**Materials Project Jain et al. [2013]:** The Materials Project is a comprehensive dataset containing material properties and structures, facilitating hypothesis generation and validation in materials science. It supports the creation of hypotheses related to material properties and applications, validates predictive models in material discovery and design, and enables structure-property relationship analysis. Evaluation metrics include structural matching to validate predicted structures, novelty filtering to identify unique materials, and prediction accuracy to assess model reliability. This dataset fosters the discovery of novel materials through computational insights and provides experimental data to validate hypotheses, ensuring practical relevance.

**KEGG Pathway Kanehisa and Goto [2000]:** The KEGG Pathway dataset offers detailed insights into metabolic and signaling pathways, supporting hypothesis generation and validation in genomics and biomedicine. It facilitates the creation of hypotheses regarding biochemical interactions, validates predictive models for gene and protein interactions, and serves as a foundation for pathway enrichment analysis. With 540 curated pathways across metabolic, regulatory, and signaling processes, the dataset is regularly updated to reflect experimental and computational advances. Evaluation metrics include pathway enrichment scores, graph metrics (centrality, connectivity, modularity), and prediction accuracy. KEGG Pathway supports the discovery of novel biochemical relationships while providing a robust basis for hypothesis validation.

**American Community Survey (ACS) Bureau [2025]:** The American Community Survey (ACS) is an annual survey by the U.S. Census Bureau, providing comprehensive demographic, social, economic, and housing data. It facilitates hypothesis generation about trends, behaviors, and disparities, validates models for policy analysis and urban development, and offers longitudinal insights into population changes. Covering approximately 2.5 million households annually and featuring over 35,000 variables, ACS enables nationwide and local-level analysis. Metrics include statistical analysis, correlation coefficients for variable relationships, and subgroup comparisons. While not suited for novel hypothesis creation, the dataset's representative and detailed nature ensures robust validation.

**Patent Data (USPTO) Patent and Office [2025]:** The USPTO dataset contains extensive information on granted patents, offering a vital resource for hypothesis generation and validation in technology and innovation research. It facilitates the exploration of technological trends, patent networks, and collaboration patterns while validating hypotheses regarding novelty and impact. With over 10 million patents spanning industries, the dataset provides metadata like inventor details, filing dates, and citations. Key metrics include citation analysis (impact scores), novelty scores, and collaboration indices. The dataset excels in novelty analysis and trend identification but is less suited for direct feasibility testing.

**XSum Narayan et al. [2018]:** The XSum dataset is a large-scale resource for abstractive text summarization, extensively used in NLP research. It supports hypothesis creation about summarization techniques, validates text generation frameworks, and benchmarks summarization algorithms. For example, a researcher might hypothesize that "pre-trained language models fine-tuned with reinforcement learning from human feedback (RLHF) produce more factually consistent summaries on XSum than models trained with maximum likelihood alone." This hypothesis can be evaluated

using XSum's single-sentence human-written summaries and news articles across diverse topics. Comprising 226,711 news articles with single-sentence summaries across diverse topics, XSum offers rich contextual diversity. Metrics include BLEU for n-gram overlap, ROUGE for unigram and sequence comparisons, and conciseness metrics. The dataset fosters novel summarization methods and provides a robust foundation for validating advanced NLP models.

**Cosmic (Catalogue of Somatic Mutations in Cancer) Bamford et al. [2004]:** Cosmic is a comprehensive dataset detailing somatic mutations in cancer, supporting hypothesis generation and validation in cancer genomics and personalized medicine. It facilitates the exploration of genetic mutations linked to cancer progression, validates mutation prediction models, and serves as a foundation for studying mutation impacts across various cancer types. The dataset includes data on over 30,000 genes and 2 million somatic mutations, annotated with pathways, phenotypes, and clinical outcomes. Evaluation metrics such as mutation analysis metrics, pathway mapping metrics, and prediction accuracy validate hypotheses. While Cosmic excels in supporting hypothesis validation, it primarily focuses on existing hypotheses and provides robust, experimentally validated data for analysis.

**Open Research Knowledge Graph (ORKG) Jaradeh et al. [2019]:** ORKG provides a structured, machine-readable representation of interdisciplinary research contributions, enabling hypothesis generation and validation. It uncovers relationships among research topics, datasets, and methods, validates hypotheses through graph-based analyses of citation impact and knowledge diffusion, and structures scientific knowledge for collaborative exploration. The dataset spans multiple disciplines, including AI, biology, and social sciences, with over 3 million triples linking research entities. Metrics such as graph centrality, novelty detection, and citation impact highlight its value in supporting novel connections and systematic validation.

**Conclusion.** The datasets presented in this section underscore their foundational role in hypothesis creation and validation across a wide range of scientific disciplines. Domain-specific datasets, such as *PubMed Abstracts* and *KEGG Pathway*, facilitate targeted research by providing structured knowledge for biomedical and genomic studies. *Interdisciplinary datasets*, like the *Open Graph Benchmark (OGB)* and *Open Research Knowledge Graph (ORKG)*, enable hypothesis generation and validation across multiple domains, fostering cross-disciplinary innovation. Datasets emphasizing novelty and feasibility, such as *Materials Project* and *ChEMBL*, support scientific breakthroughs in materials science and drug discovery by offering rich, structured data for predictive modeling and validation. Resources like *ClimateNet* and *COCO Dataset* demonstrate how curated datasets drive advancements in climate science and AI-driven hypothesis evaluation. As these datasets continue to expand in scale, annotation quality, and accessibility, they will play an increasingly vital role in enhancing hypothesis generation, ensuring rigorous validation, and accelerating scientific discovery. By leveraging these diverse and evolving datasets, researchers can refine predictive models, validate novel hypotheses, and drive breakthroughs in emerging fields.

## 4    Categorization of Hypothesis Generation Approaches

Scientific hypothesis generation has been modeled through two primary paradigms: (1) computational frameworks for discovery grounded in symbolic reasoning and cognitive science, and (2) contemporary methods driven by large-scale neural models, particularly Large Language Models (LLMs). Each reflects a different intuition about how hypotheses are conceived, represented, and evaluated. Computational frameworks for discovery view scientific hypothesis formation as a structured problem-solving process. Influenced by early work in cognitive science, these systems simulate how humans incrementally build explanations from observations by applying heuristic rules, constructing symbolic representations, and iteratively refining their models Bradshaw et al. [1983], Simon [1992], Langley and Jones [1988], Langley [1998, 2000], Džeroski et al., Langley and Simon [2013], Langley [2024]. Tools such as *BACON* Bradshaw et al. [1983] and *KEKADA* Langley [2000] exemplify this approach by rediscovering known laws and relationships in structured datasets. These methods define a hypothesis space $\mathcal{H}$ as a set of symbolic expressions—such as Newton's second law or Mendelian inheritance rules—generated through grammars, algebraic forms, or logic-based templates, often constrained by background knowledge or domain-specific primitives. Candidate hypotheses within this space are constructed using heuristic search strategies such as forward chaining Langley [2024], rule induction, or equation synthesis. They are ranked based on their empirical fit to observed data and structural simplicity, typically favoring parsimonious and generalizable formulations. This process is

often formalized as a weighted scoring function:

$$\text{score}(h) = \alpha \cdot \text{fit}(h, D) - \beta \cdot \text{complexity}(h), \qquad (7)$$

where $h \in \mathcal{H}$ is a hypothesis, $D$ is the dataset, and $\alpha$, $\beta$ are weights balancing empirical accuracy and parsimony. To illustrate, consider the task of rediscovering the Hall-Petch relationship in materials science, which relates the yield strength $\sigma_y$ of a polycrystalline material to its grain size $d$ through the equation $\sigma_y = \sigma_0 + k \cdot d^{-1/2}$. A candidate hypothesis in this domain may take the form: $h(d) = a + b \cdot d^{-c}$, where $a$, $b$, and $c$ are free parameters to be estimated from data. The empirical fit can be computed using mean squared error: $\text{fit}(h, D) = -\frac{1}{N} \sum_{i=1}^{N} \left(h(d_i) - \sigma_{y_i}\right)^2$, where $(d_i, \sigma_{y_i}) \in D$ are observed grain size and yield strength pairs. The complexity of $h$ can be quantified by counting the number of mathematical operators and the depth of the expression tree: $\text{complexity}(h) = \text{NumOperators}(h) + \lambda \cdot \text{TreeDepth}(h)$. This formalism favors hypotheses that not only fit the data well but are also structurally simple and generalizable, reflecting core principles of scientific discovery. These approaches are particularly valuable when working with well-structured, interpretable data and when the hypothesis space is tightly constrained by theory. However, their reliance on handcrafted rules and domain-specific encodings limits their scalability and effectiveness in data-rich, unstructured, or ambiguous environments.

While symbolic frameworks provide a principled and interpretable foundation for modeling hypothesis generation, they often rely on domain-specific encodings, constrained rule spaces, and hand-crafted heuristics. As scientific data's scale, heterogeneity, and ambiguity have increased, these limitations have spurred the adoption of more flexible and data-driven approaches. This shift has been catalyzed by the emergence of Large Language Models (LLMs), which enable hypothesis generation by synthesizing knowledge across unstructured sources at scale. Unlike symbolic methods that construct hypotheses from explicitly defined rule spaces, LLMs operate over implicit probabilistic representations learned from diverse corpora, offering new possibilities for discovery in underexplored or interdisciplinary domains. These models are trained on massive, heterogeneous corpora that span scientific literature, code repositories, and multimodal sources. Rather than relying on symbolic reasoning, LLMs use statistical learning to approximate the conditional probability of generating a hypothesis $H$ given a context $C$, typically modeled as:

$$P(h \mid c) = \prod_{t=1}^{T} P(h_t \mid h_{<t}, c; \theta), \qquad (8)$$

where $h = (h_1, \ldots, h_T)$ is a sequence of tokens representing a candidate hypothesis, and $\theta$ are the model parameters learned from data. LLMs excel at generalizing across domains, synthesizing knowledge from unstructured input, and proposing hypotheses that may span disciplinary boundaries. They are particularly effective in cases where structured models are not readily available or where rapid exploration of diverse ideas is desired. However, LLMs often lack formal mechanisms for explanation, causality, and logical rigor, necessitating downstream validation via simulation, symbolic reasoning, or expert review Wang et al. [2024b].

The choice between symbolic and LLM-based approaches depends largely on the nature of the task and the structure of available data. Symbolic frameworks are most effective when the objective is to derive interpretable models from well-defined variables or to formulate hypotheses aligned with established scientific theories. In contrast, LLMs are well-suited for contexts involving large-scale, unstructured datasets or where the discovery of novel, cross-domain associations is prioritized. Increasingly, hybrid pipelines that combine LLM-driven generation with symbolic or simulation-based validation are being adopted to balance generative flexibility with interpretability and rigor Langley [2024], Ghafarollahi and Buehler [2024a], Ren et al. [2025]. Recent advancements exemplify this convergence: the AHTech platform Lin et al. [2025] integrates automated electrochemical experimentation with machine learning to enable high-throughput hypothesis testing in battery research; LLMs have been shown to automate bioinformatics workflows when paired with structured repositories like cBioPortal Ji et al. [2024]; and modular systems modeling frameworks such as Robotics-LLM Yin et al. [2025] and proteomics-based KDD pipelines Resell et al. [2025] facilitate hypothesis refinement in chemical discovery and oncology, respectively. These efforts underscore modern hypothesis generation evolving into a multifaceted, interdisciplinary process where automation, experimentation, and semantic reasoning coalesce to accelerate discovery. They also highlight the growing need for domain-specific, interpretable, and safety-aware AI agents Ren et al. [2025], Yu [2025], Steinecker et al. [2025].

Table 4: Comparison of Symbolic and LLM-Based Hypothesis Generation Approaches

| Aspect | Symbolic Discovery Systems | LLM-Based Generative Systems |
|---|---|---|
| Example Systems | BACON Bradshaw et al. [1983], KEKADA Langley [2000] | ChatGPT Achiam et al. [2023], Sci-Agents Ghafarollahi and Buehler [2024b] |
| Hypothesis Space Construction | Explicit, rule-based; defined using symbolic grammars, logic rules, and algebraic templates | Implicit, data-driven; encoded via pretraining on large corpora and fine-tuning for specific domains |
| Inference Mechanism | Heuristic or search-based (e.g., forward/backward chaining, ILP) | Generative decoding using probabilistic token prediction and attention-based reasoning |
| Validation Strategy | Fit to empirical data, parsimony (Occam's razor), consistency with domain theory | Rationale scoring, entailment verification, retrieval-augmented consistency checking |
| Interpretability | High (transparent rule structures, derivations can be traced) | Low to moderate (requires explainability tools such as SHAP Lundberg and Lee [2017], LIME Ribeiro et al. [2016], or prompt engineering) |
| Scalability | Limited by combinatorial rule space and symbolic inference complexity | High scalability due to pretraining and fine-tuning across large-scale unstructured datasets |
| Typical Application Domains | Classical scientific discovery (e.g., physics, chemistry, cognitive modeling) | Interdisciplinary science, biomedical literature mining, materials science, automated experimentation |
| Strengths | Theory-grounded, interpretable, robust in structured domains | Flexible, cross-domain generalization, effective in handling unstructured or sparse data |
| Limitations | Labor-intensive to construct, brittle with noisy data, domain-dependent | Prone to hallucination, limited interpretability, sensitive to prompt and data biases |

Building on these foundations, the remainder of this section categorizes contemporary hypothesis generation approaches into distinct methodologies. These include knowledge-driven methods, data-driven integration, AI-driven exploration, text and concept mining, simulation and modeling, interactive and collaborative systems, causal inference, dynamic and adaptive knowledge systems, and multi-agent systems. Each approach reflects a unique computational intuition and contributes to the evolving landscape of scientific discovery through its specialized techniques and domain applications. By showcasing these methodologies, we aim to highlight their transformative potential in fostering novel and impactful scientific discoveries. Figure 4 presents a hypothesis generation pipeline highlighting various methods that contribute to producing candidate hypotheses that are novel and plausible.

## 4.1 Knowledge-Driven Methods

Knowledge-driven methods, which include knowledge graphs, network-based approaches, and ontology-based reasoning, represent a cornerstone of modern hypothesis generation. These methods provide a systematic and structured way to explore scientific knowledge's vast and often overwhelming complexity. Organizing information into structured representations enables researchers to uncover hidden patterns, establish interdisciplinary connections, and generate innovative and contextually relevant hypotheses. These approaches are particularly valuable in domains where the complexity of interactions—such as between genes, proteins, diseases, or materials—defies traditional manual analysis.

Knowledge graphs serve as powerful tools for capturing and visualizing relationships between entities. Knowledge graphs facilitate intuitive exploration and semantic reasoning by structuring knowledge as nodes (representing entities such as genes, diseases, or chemical compounds) and edges (representing their relationships). Systems like MOLIERE leverage vast biomedical repositories such as PubMed to identify novel gene-disease associations that often elude conventional analysis Sybrandt et al.
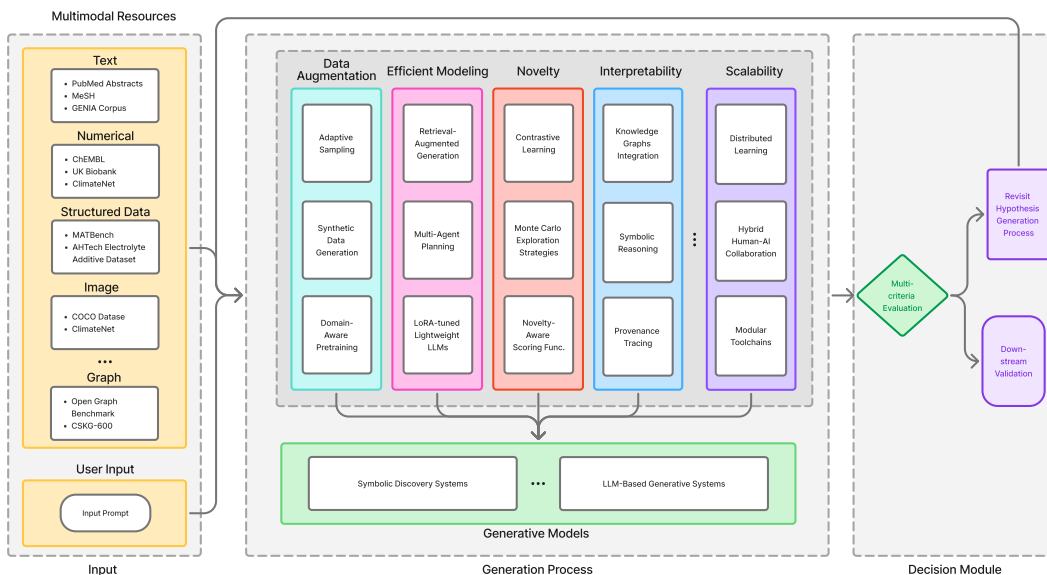
Figure 3: Modular pipeline for AI-driven hypothesis generation. The figure illustrates how multimodal data sources flow through symbolic and LLM-based generative components, incorporating retrieval, reasoning, scoring, and refinement to support interpretable and novelty-aware hypothesis generation.

[2018]. Similarly, SciAgents demonstrates the integration of dynamic knowledge graphs with large language models (LLMs), enabling interdisciplinary research in fields like pharmacology, where complex relationships must be navigated to propose innovative hypotheses Ghafarollahi and Buehler [2024b]. These systems allow researchers to uncover overlooked connections, prioritize promising research directions, and bridge gaps across disciplines, accelerating discovery. Recent work has also explored the reverse paradigm: using LLMs themselves as latent knowledge graphs. Instead of explicitly constructing graph structures, these approaches treat LLMs as implicit knowledge stores capable of retrieving, composing, and reasoning over relational information. For instance, Kau et al. [2024] proposes prompting strategies and architectural mechanisms that enable LLMs to emulate knowledge graph behavior by performing multi-hop reasoning or generating structured triples directly from natural language inputs. While this direction offers scalability and ease of deployment advantages, reducing the need for manual graph construction, it often sacrifices interpretability, semantic precision, and logical consistency. LLM-based representations lack the explicit structure and auditability of formal knowledge graphs, making them less suitable for domains requiring rigorous semantic alignment. Therefore, these approaches are best viewed as complementary; LLMs can enhance knowledge graph population, reasoning, and hypothesis suggestion, but structured graphs remain crucial for precise and semantically grounded scientific hypothesis generation.

A key strength of knowledge graphs lies in their ability to perform advanced analyses such as link prediction and graph-based learning. Link prediction, which leverages metrics like embeddings, centrality, and similarity scores, has proven particularly impactful in areas like drug discovery, where identifying new drug-disease interactions can expedite therapeutic development Kim and Segev [2018]. Recent advancements, such as Graph Neural Networks (GNNs), further extend the power of knowledge graphs by modeling intricate, high-dimensional relationships Bai et al. [2024a]. GNNs excel in revealing latent patterns and dependencies, enabling the analysis of multi-layered interactions that drive innovation in fields like materials science, environmental modeling, and beyond Beltagy et al. [2019]. Furthermore, the dynamic nature of modern knowledge graphs ensures their adaptability, as they can incorporate newly available data to refine their structure and maintain relevance in rapidly evolving scientific landscapes Qi et al. [2024].

Ontology-based reasoning complements knowledge graphs by introducing formalized frameworks for semantic consistency and logical reasoning. Unlike knowledge graphs, which focus on uncovering relationships, ontology-based systems emphasize ensuring that generated hypotheses adhere to established domain-specific standards and terminologies. These systems use well-defined ontologies to capture classes, properties, and relationships, thereby enabling precise and semantically valid

Table 5: Summary of Hypothesis Generation Approaches and Tools

| Approach | Example Tool/System | Domain | Strengths | Weaknesses |
|---|---|---|---|---|
| Knowledge-Driven Methods | MOLIERE Sybrandt et al. [2018] | Biomedicine, Interdisciplinary | Ensures consistency and logical validity | Limited novelty, relies on existing knowledge |
| Data-Driven Integration | SciAgents Ghafarollahi and Buehler [2024b] | Biomedicine, Scientific Discovery | Merges structured and unstructured data | Quality depends on dataset reliability |
| AI-Driven Exploration | Reinforcement Learning Zhou et al. [2024] | Drug Discovery, Materials Science | Uncovers novel patterns beyond human intuition | High computational cost, lacks interpretability |
| Text and Concept Mining | Chemist-X Chen et al. [2024] | Chemistry, Biomedical Sciences | Extracts insights from large-scale text data | Sensitive to NLP noise, domain adaptation required |
| Simulation and Modeling | VELMA Schumann et al. [2024] | Robotics, Engineering | Models diverse scenarios and unconventional ideas | Requires high computational resources |
| Interactive and Collaborative Systems | Human-AI Collaboration Kim and Segev [2018] | Education, Ethics | Combines human insights with AI efficiency | Dependent on human oversight, potential biases |
| Causal Inference | Causal Discovery Jha et al. [2019] | Biomedicine, Social Science | Identifies causal relationships, enhances reliability | Requires large datasets, vulnerable to confounders |
| Dynamic and Adaptive Knowledge Systems | Knowledge Graph Updates Beltagy et al. [2019] | Pharmacology, AI Research | Continuously refines and contextualizes knowledge | High processing power needed for real-time updates |
| Multi-Agent Systems | Multi-Agent AI Qi et al. [2024] | Genomics, Collaborative Science | Enables decentralized, expert-driven collaboration | Complexity increases with agent interactions |

hypothesis generation. For example, ontology-based hypothesis validation has been used to identify novel gene-disease relationships by mapping genetic functions to clinical phenotypes. The integration of multiple ontologies has further unified knowledge across domains such as physics, biology, and chemistry, facilitating interdisciplinary research Wang et al. [2023].

The methodologies underpinning ontology-based reasoning include ontology mapping, semantic inference, and automated reasoning. Ontology mapping aligns diverse ontologies from multiple domains, enabling unified views of knowledge essential for interdisciplinary research Kim and Segev [2018]. Semantic inference, which applies logical rules to infer new relationships, uncovers insights that are otherwise difficult to detect. Automated reasoning tools derive novel hypotheses while maintaining consistency with existing knowledge structures Beltagy et al. [2019]. These techniques have demonstrated significant utility in biomedical science Ji et al. [2024], drug discovery Blanco-Gonzalez et al. [2023], and interdisciplinary innovation Zhang et al. [2024b], with applications ranging from identifying novel drug targets to exploring complex ecological systems Qi et al. [2024].

Together, knowledge graphs and ontology-based reasoning form a complementary toolkit for hypothesis generation. Knowledge graphs Bai et al. [2024b] excel at uncovering patterns and relationships that are often hidden in vast, unstructured datasets, while ontology-based reasoning Liu et al. [2010] ensures that hypotheses are both precise and semantically consistent. This synergy has profound implications across disciplines. In biomedicine, these methods have mapped intricate relationships between genes, diseases, and drugs, facilitating breakthroughs in personalized medicine and rare disease research Kim and Segev [2018]. In materials science, they have accelerated the discovery of novel materials by identifying promising chemical combinations Ghafarollahi and Buehler [2024a].

Environmental science has also benefited from these approaches, with knowledge graphs uncovering correlations between oceanic and atmospheric variables to model ecological phenomena Wang et al. [2023]. Ontology-based reasoning has further strengthened interdisciplinary research by aligning diverse terminologies and frameworks, enabling seamless collaboration and hypothesis generation across fields.

While knowledge-driven methods have proven transformative, they are not without limitations. Knowledge graphs rely heavily on structured data, which can constrain their ability to generate truly novel hypotheses when data is incomplete or biased. Similarly, ontology-based reasoning is limited by the granularity and completeness of the underlying ontologies. To address these challenges, researchers are increasingly integrating dynamic updates and machine learning techniques. For instance, GNNs and other machine learning models enhance the ability of knowledge graphs to explore higher-order interactions and adapt to new data. Ontology-based systems are also evolving, incorporating learning-based methods to identify less obvious relationships and expand their scope. These advancements ensure that knowledge-driven methods remain adaptable, innovative, and capable of addressing the complexities of modern scientific inquiry.

In conclusion, knowledge-driven methods are indispensable for navigating and extracting value from the vast complexity of scientific knowledge. By combining structured representations, semantic rigor, and dynamic adaptability, they provide researchers with the tools to uncover hidden patterns, bridge disciplines, and generate actionable hypotheses. As these methods continue to evolve, their potential to drive innovation and foster interdisciplinary collaboration remains boundless.

## 4.2 Data-Driven Integration

Data-driven integration represents a transformative approach to hypothesis generation by synthesizing diverse datasets from disciplines such as genomics, proteomics, environmental science, and beyond. This methodology bridges traditionally isolated domains, uncovering complex relationships and enabling researchers to address multifaceted scientific challenges with unprecedented precision. By leveraging cross-domain data integration, researchers can construct a holistic understanding of intricate systems, fostering innovative solutions and advancing scientific discovery across a wide array of fields.

At the heart of data-driven integration lie methodologies designed to analyze and synthesize heterogeneous datasets, revealing insights that are otherwise inaccessible. Multi-omics integration, for instance, combines data from genomics, transcriptomics, and proteomics to identify high-value targets such as gene-protein interactions Qi et al. [2024]. This approach has driven breakthroughs in personalized medicine by uncovering molecular mechanisms underlying diseases and identifying actionable therapeutic targets. Beyond omics, interdisciplinary integration extends these principles to link disparate datasets, such as combining genomic information with environmental data to explore broader patterns and emergent properties Pammi et al. [2023]. These techniques often rely on systems-level computational models, which simulate interactions across disciplines, providing a dynamic framework for hypothesis generation Kim and Segev [2018]. By capturing non-obvious correlations and emergent behaviors, these models empower researchers to explore complex phenomena that transcend traditional disciplinary boundaries.

A growing suite of tools exemplifies the far-reaching impact of data-driven integration in scientific hypothesis generation across diverse domains. In agriculture, VirtualPlant Katari et al. [2010] integrates genomic, transcriptomic, and phenotypic data to uncover genetic pathways that enhance crop resistance to environmental stressors, supporting sustainable farming practices Qi et al. [2024], Kim and Segev [2018]. In pharmacology, BioLunar Wysocki et al. [2024] leverages multi-omics data to elucidate complex gene-protein interactions, enabling the discovery of precision therapies and accelerating the identification of drug-protein interactions Kim and Segev [2018], Qi et al. [2024]. In environmental science, tools like Climate KG Wu et al. [2022] link climate variables with ecological and biological datasets to uncover factors driving ecosystem resilience and inform conservation and adaptation strategies. These examples not only demonstrate the versatility of integrative tools in enabling hypothesis generation but also highlight their broader utility in addressing critical challenges across genomics, pharmacology, agriculture, and environmental science.

The novelty of data-driven integration lies in its ability to connect datasets traditionally analyzed in isolation. For example, integrating ocean salinity data with atmospheric pressure measurements has

provided new insights into climate dynamics that were previously unattainable Wang et al. [2023]. However, this potential for discovery is contingent on the availability of high-quality datasets and the computational resources required to process and interpret them effectively. Resource constraints and dataset heterogeneity pose challenges that must be addressed to fully realize the promise of data-driven integration Gruver et al. [2024]. To overcome these barriers, adaptive algorithms, reinforcement learning, and dynamic systems modeling are increasingly employed. These methods not only enhance computational efficiency but also enable the exploration of deeper, non-obvious correlations within interdisciplinary datasets Zhou et al. [2024].

Data-driven integration serves as a powerful paradigm for addressing complex scientific questions, bridging disciplinary divides, and fostering novel insights. As computational methods and data quality continue to improve, this approach is poised to redefine the landscape of hypothesis generation. By synthesizing diverse data sources into cohesive frameworks, data-driven integration empowers researchers to tackle intricate challenges and unlock discoveries that were once beyond reach. The continued evolution of these methodologies promises to drive progress across a wide spectrum of fields, solidifying their role as a cornerstone of modern scientific inquiry.

## 4.3 AI-Driven Exploration

AI-driven exploration represents a transformative paradigm in hypothesis generation, integrating advanced methodologies such as machine learning (ML), statistical modeling, retrieval-augmented generation (RAG), and reinforcement learning (RL). These approaches empower researchers to process vast and complex datasets, synthesize real-time information, and navigate intricate hypothesis spaces to uncover innovative solutions. By leveraging adaptive algorithms, dynamic retrieval mechanisms, and iterative refinement strategies, AI-driven methods address the challenges of modern scientific inquiry across fields such as drug discovery, materials science, robotics, and social sciences Qi et al. [2024], Gruver et al. [2024].

RAG exemplifies the dynamic capabilities of AI in hypothesis generation by combining the language comprehension power of LLMs with domain-specific datasets and real-time data retrieval systems. This approach bridges the gap between static repositories and evolving research challenges, ensuring that hypotheses are timely and grounded in relevant data Zhou et al. [2024]. Key techniques in RAG include dynamic information retrieval, prompt engineering, and adaptive memory integration. Dynamic retrieval enables the extraction of up-to-date information from diverse sources, enriching hypotheses with the latest knowledge Beltagy et al. [2019]. Carefully crafted prompts guide LLMs to generate contextually specific outputs, addressing discipline-specific challenges with precision Kim and Segev [2018]. Memory integration further enhances RAG's capabilities by incorporating feedback from previous outputs, fostering coherence and novelty Qi et al. [2024]. Tools like Chemist-X and SciAgents exemplify RAG's versatility, applying it to drug discovery, interdisciplinary research, and pharmacology Chen et al. [2024], Ghafarollahi and Buehler [2024b].

RL adds another dimension to AI-driven exploration by leveraging trial-and-error strategies to refine hypotheses dynamically. RL systems optimize solutions in high-dimensional and non-linear hypothesis spaces, enabling the discovery of novel insights and the refinement of processes Fawzi et al. [2022]. Techniques such as policy optimization iteratively improve decision-making processes, while model-free RL approaches like Q-learning facilitate hypothesis generation in environments with complex or poorly understood dynamics Gruver et al. [2024]. Multi-agent RL systems expand this capability further by employing collaborative agents to explore diverse hypothesis spaces simultaneously, fostering interdisciplinary insights and significantly increasing efficiency Zhou et al. [2024].

Several tools showcase the practical applications of these methodologies. VELMA demonstrates the integration of RAG with navigation and robotics, using textual and visual data to hypothesize optimal strategies for urban navigation Schumann et al. [2024]. Chemist-X applies RAG to chemical databases, generating reaction pathways and identifying novel drug candidates Chen et al. [2024]. In RL, tools like RL-Discovery optimize material properties by exploring hypothesis spaces, accelerating advancements in catalyst design Gruver et al. [2024]. DrugRL leverages Q-learning to refine molecular configurations in drug discovery, while SciAgents integrates RL and knowledge graph approaches to enhance hypothesis validation across interdisciplinary fields Ghafarollahi and Buehler [2024b].

AI-driven exploration has demonstrated transformative impact across scientific disciplines. In drug discovery, RL frameworks such as DrugRL have accelerated the design of molecular structures, reducing costs and timelines by iteratively refining candidate molecules Qi et al. [2024]. In materials science, RL-Discovery has optimized catalysts and identified compounds with unique properties, addressing critical performance criteria Gruver et al. [2024]. RAG tools like SciAgents have facilitated interdisciplinary queries, enabling researchers to tackle biomedical and computational challenges through systematic hypothesis refinement Ghafarollahi and Buehler [2024b]. These examples highlight the adaptability and efficacy of AI-driven exploration in addressing complex, multi-dimensional problems.

While these methodologies hold immense promise, they are not without challenges. RAG systems depend on the quality and diversity of retrieved data, and poorly curated sources can constrain the novelty of hypotheses Beltagy et al. [2019]. RL systems, meanwhile, are heavily influenced by the design of reward functions, with poorly defined rewards potentially leading to suboptimal exploration Fawzi et al. [2022]. To address these limitations, researchers are integrating complementary approaches such as Bayesian optimization and adaptive memory mechanisms to balance exploration and exploitation while enhancing creativity and scientific relevance Zhou et al. [2024].

In conclusion, AI-driven exploration integrates RAG, RL, and advanced ML techniques into a cohesive framework for hypothesis generation. By navigating dynamic data landscapes, refining hypotheses iteratively, and uncovering high-value discoveries, these methodologies are expanding the frontiers of scientific inquiry. As AI-driven systems continue to evolve, their potential to foster interdisciplinary collaboration, drive innovation, and accelerate breakthroughs across diverse fields remains boundless.

## 4.4 Text and Concept Mining

Text and concept mining have emerged as indispensable methodologies for extracting meaningful insights from vast repositories of unstructured textual data, such as scientific literature, patents, and technical reports. These approaches systematically analyze and track the evolution of concepts, enabling researchers to uncover latent patterns, identify emerging trends, and generate hypotheses that align with the ever-changing knowledge landscape. By bridging information across disciplines and highlighting connections that may otherwise remain obscured, text and concept mining empower scientific discovery and innovation Pajo [2025].

At the core of text and concept mining are advanced techniques designed to distill valuable insights from massive text corpora. Topic modeling algorithms, such as Latent Dirichlet Allocation (LDA), cluster related information into coherent themes, revealing underlying patterns that might not be immediately apparent Beltagy et al. [2019]. Dynamic evolution models extend these capabilities by analyzing temporal changes in text representations, such as word embeddings or semantic networks, to capture shifts in research priorities and conceptual frameworks Qi et al. [2024]. These longitudinal analyses are critical for identifying emerging fields of study and generating hypotheses that reflect current and future trends. Natural Language Processing (NLP) pipelines further enhance the process by automating tasks like text preprocessing, feature extraction, and semantic analysis, ensuring that the generated hypotheses are both contextually relevant and actionable Zhou et al. [2024].

A range of innovative tools exemplifies the transformative potential of text and concept mining in hypothesis generation. Dyport combines text mining with dynamic graph modeling to uncover temporal trends and co-occurrence networks in genomics, enabling the identification of novel gene-disease associations Tyagin and Safro [2024]. SciFact leverages co-occurrence patterns in scientific literature to propose new drug-disease relationships, opening avenues for drug repurposing. ConceptNet is a pre-built common-sense knowledge base, helping researchers identify knowledge gaps and propose innovative hypotheses in diverse fields Kim and Segev [2018]. These tools collectively demonstrate how text and concept mining methodologies can address complex scientific challenges by providing a dynamic framework for discovery.

The applications of text and concept mining span a wide array of scientific disciplines, illustrating their versatility and impact. In genomics, Dyport has analyzed temporal trends in gene-variant literature, revealing emerging links to rare diseases and prioritizing research directions Tyagin and Safro [2024]. In drug discovery, SciFact has enabled the rapid identification of novel drug-disease pairings through large-scale analysis of biomedical texts, significantly accelerating the exploration

of repurposing opportunities. In climate science, dynamic evolution models have tracked decades of literature to highlight shifts in research priorities, methodologies, and key focus areas, offering critical insights into the evolution of the field. These examples underscore the adaptability of text and concept mining in addressing interdisciplinary and domain-specific challenges.

The novelty of hypotheses generated by text and concept mining stems from their ability to uncover emerging connections and align with cutting-edge research directions. By analyzing temporal trends and co-occurrence patterns, these methods reveal hidden insights and suggest promising avenues for exploration. However, their effectiveness depends on the quality, recency, and completeness of the underlying datasets. Outdated or incomplete text corpora can constrain the scope and relevance of generated hypotheses, limiting their potential impact. To overcome these limitations, advancements such as multilingual and multimodal text mining, which combine textual data with other modalities such as images, tables, or graphs to uncover richer insights, are being integrated, broadening the diversity of insights and fostering cross-disciplinary hypothesis generation. These enhancements ensure that text and concept mining remain robust tools for driving scientific innovation.

By harnessing the vast and ever-expanding repositories of textual data, text and concept mining provide researchers with a powerful, adaptive framework for hypothesis generation. Their ability to track the evolution of knowledge, uncover emerging trends, and propose novel connections positions them as essential methodologies in the modern scientific toolkit. As these techniques continue to evolve, their role in advancing research across disciplines is poised to expand, catalyzing progress and innovation in diverse scientific domains.

## 4.5 Simulation and Modeling

Simulation and modeling approaches Schumann et al. [2024], Lavin et al. [2021], Sarathy and Scheutz [2022], Clark et al. [2020], Hélie and Sun [2008] are revolutionizing hypothesis generation by mimicking creative problem-solving processes to explore unconventional ideas and solutions. These techniques leverage advanced computational frameworks to emulate divergent thinking, allowing researchers to transcend traditional methodologies and tackle complex challenges with fresh perspectives. By enabling exploration across vast solution spaces, simulation and modeling foster innovation and uncover opportunities that might otherwise remain hidden. This approach has proven particularly effective in addressing non-linear, multi-dimensional problems across diverse scientific disciplines Kim and Segev [2018], Qi et al. [2024].

Central to simulation and modeling are techniques designed to simulate creativity and enhance exploratory reasoning. Simulated annealing, inspired by natural processes of energy optimization, iteratively balances exploration and exploitation to identify optimal hypotheses. This method ensures thorough exploration of solution spaces while converging toward high-value outcomes. Creative neural networks, leveraging specialized architectures, emulate human-like reasoning and imagination to propose novel and unconventional hypotheses. These networks excel in synthesizing innovative ideas, making them particularly valuable in domains where groundbreaking insights are essential. Idea mining, which combines text mining and clustering techniques, extracts unique concepts from large datasets, enhancing the diversity and originality of generated hypotheses. Together, these methodologies provide a powerful framework for generating novel and impactful ideas.

A variety of cutting-edge tools demonstrate the transformative potential and practical versatility of simulation and modeling in hypothesis generation across diverse domains. In the social sciences, SimHypoth Lin and Lucas [2023] employs neural network-based models to uncover non-linear relationships and generate hypotheses that challenge conventional assumptions, including correlations between demographic factors and societal trends. In engineering and materials science, IdeaFlow Utley and Klebahn [2022] tracks the diversity of ideas emerging from simulations, enabling the discovery of innovative design strategies and material combinations that have advanced sustainable infrastructure, including novel approaches to bridge and building construction. In theoretical physics, GenCreative de Oliveira et al. [2017] leverages generative adversarial networks (GANs) to explore unconventional solutions to unsolved problems, from subatomic particle interactions to alternative theoretical frameworks, thereby expanding both theoretical inquiry and experimental frontiers. Collectively, these tools exemplify how simulation and modeling techniques can drive scientific innovation and address complex, interdisciplinary challenges.

The novelty of hypotheses generated through simulation and modeling is one of their defining strengths. By embracing unconventional approaches, these techniques frequently lead to insights that challenge existing paradigms and expand the boundaries of current knowledge. For instance, IdeaFlow proposed novel material combinations in construction that had not been previously considered, demonstrating the power of creative exploration Kim and Segev [2018]. However, the reliance on randomness or unstructured exploration in some simulations can occasionally produce impractical or untestable hypotheses, limiting their utility. To mitigate this, researchers increasingly incorporate domain-specific constraints and iterative feedback loops, ensuring that hypotheses remain both innovative and feasible. These enhancements balance creativity with relevance, allowing simulation and modeling to generate actionable insights.

Simulation and modeling represent a dynamic frontier in hypothesis generation, blending computational power with creativity to tackle complex scientific challenges. By fostering unconventional thinking and expanding solution spaces, these approaches provide a robust framework for advancing scientific inquiry and driving innovation across disciplines. As techniques and tools continue to evolve, their capacity to generate transformative insights will remain central to addressing the complexities of modern research.

## 4.6 Interactive and Collaborative Systems

Interactive and collaborative systems, including human-in-the-loop (HITL) frameworks, combine the intuitive reasoning of human experts with the computational power of AI, creating a synergistic approach to hypothesis generation. By integrating iterative feedback from users, these systems refine AI-generated outputs to ensure hypotheses are relevant, innovative, and aligned with real-world contexts. This blend of human creativity and machine efficiency enables researchers to explore complex problems, bridging gaps between intuition-driven insights and data-driven precision Kim and Segev [2018], Ghafarollahi and Buehler [2024b].

The methodologies driving interactive and collaborative systems focus on iterative learning, crowd-sourced contributions, and explainability. Iterative learning fosters a continuous feedback loop between human users and AI, allowing hypotheses to be refined through multiple cycles of interaction. This ensures that the outputs align with the strategic priorities and nuanced understanding of domain experts Zhou et al. [2024]. Crowdsourced contributions expand this collaborative framework by aggregating input from diverse participants, enriching the hypothesis space with multiple perspectives and enhancing robustness Kim and Segev [2018]. Explainability mechanisms ensure that AI-generated hypotheses are interpretable, providing transparent justifications for their derivation. This fosters trust and enables meaningful feedback, ultimately improving the quality and credibility of the generated hypotheses Beltagy et al. [2019].

A variety of tools exemplify the practical implementation of these methodologies, showcasing their versatility across scientific domains. SciAgents integrates real-time expert feedback with AI-driven hypothesis generation, enabling the refinement of complex hypotheses in biomedical and pharmacological research Ghafarollahi and Buehler [2024b]. CrowdScience leverages crowdsourcing to generate hypotheses on behavioral and social phenomena, aggregating insights from a wide participant pool to ensure inclusivity and diversity. ExplanatoryAI focuses on interpretable AI-driven hypothesis generation, ensuring that hypotheses are both transparent and actionable, making it particularly valuable in ethically sensitive fields such as governance and policy-making Shavit et al. [2023]. These tools demonstrate the power of interactive and collaborative systems to drive innovation by combining human expertise with AI-driven exploration.

The practical applications of interactive and collaborative systems span a wide array of disciplines, underscoring their transformative potential. In biomedical research, SciAgents has accelerated the identification of disease biomarkers and drug-target interactions by integrating expert feedback, contributing to faster development of diagnostic and therapeutic solutions Ghafarollahi and Buehler [2024b]. In social sciences, CrowdScience has facilitated the collaborative generation of hypotheses on behavioral trends and social correlations, incorporating diverse stakeholder insights to uncover novel patterns Zhou et al. [2024]. In ethics and governance, ExplanatoryAI has advanced policy development and regulatory frameworks by ensuring that hypotheses are interpretable and actionable, aiding decision-makers in evaluating the potential impacts of their choices Shavit et al. [2023]. These examples highlight the adaptability of interactive and collaborative systems in addressing domain-specific and interdisciplinary challenges.

The novelty of hypotheses generated by these systems lies in their ability to integrate human insight with AI-driven exploration. Human expertise provides contextual understanding, while AI expands the hypothesis space, enabling the discovery of innovative and practical solutions Zhou et al. [2024]. However, reliance on human input can introduce biases or constrain the exploration of unconventional ideas Kim and Segev [2018]. To address this, recent advancements incorporate generative adversarial techniques (GANs) to balance human feedback with machine-generated alternatives, fostering a broader and more creative exploration of hypotheses Beltagy et al. [2019]. These enhancements ensure that hypotheses are not only relevant but also push the boundaries of conventional thinking.

Interactive and collaborative systems represent a powerful paradigm for hypothesis generation, blending the creativity and expertise of humans with the computational efficiency and scalability of AI. By enabling iterative collaboration, fostering inclusivity, and ensuring transparency, these systems drive the creation of actionable and innovative hypotheses, paving the way for breakthroughs across diverse scientific fields.

## 4.7 Causal Inference

Causal inference frameworks Khatibi et al. [2024], Peters et al. [2017], Lucas [2007], Neuberg [2003] leverage advanced statistical and computational techniques to identify and analyze causal relationships within datasets. Unlike correlation-based approaches, which often fail to capture the underlying mechanisms of observed phenomena, causal inference frameworks focus on uncovering cause-and-effect pathways. This makes them particularly effective for generating actionable hypotheses grounded in causality, enabling researchers to propose interventions and predict their outcomes with greater confidence. By bridging the gap between observational data and experimental insights, causal inference frameworks have become indispensable tools across diverse scientific domains Jha et al. [2019], Qi et al. [2024].

At the core of causal inference methodologies are techniques designed to map, test, and validate causal relationships. Structural Causal Models (SCMs) represent one of the foundational approaches, using directed acyclic graphs (DAGs) to visualize and analyze causal pathways Jha et al. [2019]. These models allow researchers to identify direct and indirect relationships between variables, providing a comprehensive framework for understanding complex systems. Bayesian networks complement SCMs by incorporating probabilistic reasoning to estimate the likelihood of causal links between variables Kim and Segev [2018]. This approach is particularly valuable in scenarios with uncertainty or incomplete information, where Bayesian reasoning helps prioritize plausible hypotheses. Interventional analysis takes causal inference a step further by simulating potential interventions within datasets to test causal hypotheses and predict the outcomes of specific actions Ghafarollahi and Buehler [2024b]. Together, these techniques form a robust foundation for hypothesis generation, enabling researchers to derive mechanistic insights from observational data.

Several cutting-edge tools demonstrate the versatility and transformative potential of causal inference frameworks in hypothesis generation across diverse scientific domains. In biomedical research, CausalNet Peters et al. [2017], which is based on Structural Causal Models (SCMs), maps complex causal pathways to generate hypotheses about disease progression and treatment effects, facilitating advances in understanding disease mechanisms and enabling predictive diagnostics and targeted therapies. In social and biomedical sciences, BayesCausality Lucas [2007] leverages Bayesian networks to uncover causal relationships between socio-economic variables and health outcomes, providing robust, evidence-based insights for public health and policy interventions. In materials science, InterveneAI Neuberg [2003] uses interventional analysis to simulate experimental conditions and identify causal relationships between material properties and performance metrics, accelerating the design and optimization of high-performance materials. Together, these tools illustrate how causal inference frameworks can generate mechanistically grounded hypotheses and support data-driven decision-making across biomedicine, social science, and engineering.

The novelty of hypotheses generated by causal inference frameworks lies in their ability to uncover mechanistic insights often overlooked by traditional methods. By focusing on causal relationships, these frameworks have identified unexpected treatment effects in biomedical datasets, challenging existing paradigms and opening new avenues for research Jha et al. [2019]. However, the effectiveness of these frameworks is contingent on the quality of the input data. Incomplete datasets or the presence of latent variables can limit their capacity to generate innovative hypotheses, underscoring the importance of robust data collection and preprocessing Shavit et al. [2023]. To enhance novelty,

researchers are increasingly integrating causal inference with machine learning techniques, such as deep generative models. This combination allows for the simultaneous capture of causal and non-causal patterns, broadening the scope of hypothesis generation while maintaining scientific rigor Beltagy et al. [2019].

Causal inference frameworks provide a powerful approach to understanding the underlying mechanisms of complex systems. By combining robust methodologies, advanced tools, and interdisciplinary applications, these frameworks enable researchers to generate actionable and novel hypotheses that drive scientific discovery and innovation. As the integration of causal inference with machine learning continues to evolve, the potential for these frameworks to redefine research paradigms and address pressing scientific challenges remains immense.

## 4.8    Dynamic and Adaptive Knowledge Systems

Dynamic and adaptive knowledge systems Fecho et al. [2021], Yan and Chen, Mosbach et al. [2020], Rossi et al. [2020], particularly Dynamic Knowledge Graphs (DKGs), revolutionize hypothesis generation by incorporating temporal and real-time updates into traditional knowledge graph frameworks. Unlike static datasets, which can quickly become outdated, DKGs adapt to new data as it emerges, providing a flexible and responsive platform for capturing the evolving dynamics of scientific knowledge. This adaptability is critical in fast-paced fields such as biomedicine, environmental science, and chemical engineering, where maintaining the relevance and accuracy of hypotheses is paramount Kim and Segev [2018], Qi et al. [2024].

The methodologies underlying DKGs emphasize the integration of temporal, contextual, and real-time information to enhance hypothesis generation. Graph Temporal Networks (GTNs) serve as a foundational technique, encoding time-stamped edges to represent evolving relationships between entities. This enables researchers to uncover temporal patterns and shifts in data that static models might overlook. Real-time data integration dynamically updates graph structures, ensuring that hypotheses are continuously informed by the most current knowledge Gu and Krenn [2025]. This capability is particularly valuable in disciplines where rapid discoveries demand immediate incorporation of new findings. Contextual embeddings further refine hypothesis precision by adapting graph representations to reflect the changing contexts of entities and their relationships, allowing for more nuanced and accurate insights. Together, these methodologies provide DKGs with the agility to address the dynamic needs of modern scientific inquiry.

Several advanced tools demonstrate the transformative potential and interdisciplinary applicability of Dynamic Knowledge Graphs (DKGs) in scientific research. In biomedicine, UpdateKG Fecho et al. [2021] is a dynamic graph system that enables real-time tracking and refinement of gene-disease associations, supporting genomic research and personalized medicine by keeping hypotheses aligned with emerging discoveries. In environmental science, SciGraph Yan and Chen dynamically models ecological interactions, such as predator-prey dynamics and environmental conditions, to generate actionable insights into ecosystem shifts and resilience. In chemical engineering, KG-Stream Mosbach et al. [2020] integrates streaming experimental data to predict chemical properties and refine hypotheses, accelerating the design and discovery of innovative materials and industrial processes. These tools highlight the adaptability of DKGs in addressing complex challenges across domains and underscore their value in generating timely, data-informed scientific hypotheses.

The novelty of hypotheses generated through DKGs lies in their ability to adapt to evolving datasets and uncover time-sensitive patterns. By leveraging time-aware updates, DKGs capture trends and insights that static systems often miss, such as transient gene-disease associations or short-lived environmental phenomena. However, reliance on real-time data streams introduces challenges, including the potential for incomplete or delayed updates that may limit the accuracy of generated hypotheses. To mitigate these limitations, integrating DKGs with predictive modeling techniques has emerged as a powerful enhancement. Predictive models enable DKGs to anticipate future trends, broadening the scope of hypothesis generation and fostering forward-looking insights.

Dynamic and adaptive knowledge systems represent a significant advancement in the field of hypothesis generation. By seamlessly integrating temporal and contextual information with real-time updates, these systems provide researchers with a robust and flexible framework for exploring complex scientific questions. Their ability to generate hypotheses that are both innovative and actionable

makes DKGs indispensable tools in the modern scientific landscape, driving progress across a diverse array of disciplines.

## 4.9  Multi-Agent Systems

Multi-agent systems (MAS) Su et al. [2024], Baek et al. [2024], Kurakin and Bredesen [2007], Park et al. [2024] represent an innovative approach to hypothesis generation by deploying multiple autonomous agents that work collaboratively to explore complex hypothesis spaces. By assigning specialized roles to individual agents and fostering interactions among them, MAS enable the generation of hypotheses enriched by diverse expertise and cross-disciplinary insights. This collaborative dynamic mirrors the human scientific process, where teams with varied backgrounds and skills work together to address multifaceted challenges. MAS offer a scalable and efficient framework for hypothesis generation, particularly in domains requiring the integration of vast and heterogeneous datasets.

The techniques underpinning MAS focus on role specialization, collaboration, and knowledge sharing to maximize the efficiency and creativity of hypothesis generation. Role-based specialization assigns distinct responsibilities to agents based on their domain expertise, ensuring that each agent focuses on a targeted aspect of the problem. This division of labor reduces redundancy and enhances the relevance of generated hypotheses. Negotiation and collaboration among agents enable the synthesis of partial hypotheses into cohesive and comprehensive propositions. By leveraging diverse perspectives, these interactions produce hypotheses that are more robust and interdisciplinary. Distributed knowledge sharing further accelerates the discovery process by integrating multi-domain insights, allowing agents to build upon each other's findings and continuously refine their outputs. Together, these techniques make MAS a powerful tool for navigating complex scientific landscapes.

Several Multi-Agent System (MAS) frameworks exemplify the versatility and transformative potential of this methodology in hypothesis generation across scientific disciplines. In genomic research, TAIS (Team AI Scientists) Liu et al. [2024] enables agents to collaboratively analyze gene expression data and generate hypotheses on gene-protein interactions, advancing our understanding of molecular biology. In materials science, the MAS framework developed by Park et al. [2024] allows agents to explore chemical spaces and identify catalysts with unique properties, generating hypotheses that balance chemical feasibility with performance metrics and accelerating the discovery of high-performance materials. ResearchAgent Baek et al. [2024] simulates collaborative research environments in interdisciplinary settings, enabling the synthesis of insights from domains such as physics, chemistry, and biology to address complex, cross-domain scientific questions. A recent addition, VirSci Su et al. [2024], is a large-scale, LLM-based MAS designed to simulate real-world scientific collaboration. It assembles virtual scientist agents with distinct backgrounds to participate in structured discussions, novelty assessments, and iterative abstract refinement, showing significant gains in generating original research ideas compared to single-agent and prior multi-agent baselines. These frameworks collectively highlight MAS's adaptability and capacity to drive innovation in modern scientific discovery.

The novelty of hypotheses generated by MAS lies in their ability to leverage diverse expertise and uncover multi-domain relationships. By integrating knowledge from different fields, MAS can propose groundbreaking hypotheses, such as linking gene expression patterns to disease phenotypes Kim and Segev [2018]. However, the novelty of these hypotheses can be constrained by communication and coordination overhead, which may hinder scalability and reduce the system's ability to address highly complex problems Zhou et al. [2024]. To overcome these limitations, advanced negotiation protocols and enhanced knowledge-sharing mechanisms have been developed. These enhancements foster deeper interdisciplinary interactions, promoting the generation of more innovative and robust hypotheses Beltagy et al. [2019].

Multi-agent systems represent a dynamic and collaborative framework for hypothesis generation, mirroring the diverse expertise and teamwork characteristic of human research teams. By integrating advanced techniques for collaboration and knowledge sharing, MAS has the potential to revolutionize scientific discovery, providing scalable and efficient solutions for addressing the complexities of modern research.

**Conclusion.** The diverse approaches to hypothesis generation demonstrate the transformative potential of LLMs and computational techniques in advancing scientific discovery. Each method brings

unique strengths: knowledge graphs uncover conceptual relationships, text mining highlights trends, machine learning drives novelty, retrieval-augmented generation enhances relevance, and multi-omics integration reveals insights into complex systems. Despite challenges like data constraints and interpretability, these methods promise to democratize hypothesis generation and accelerate break-throughs across disciplines. Integrating these methodologies, such as combining machine learning with ontology-based reasoning or dynamic knowledge graphs with multi-agent systems, amplifies their strengths. Such hybrid strategies enable comprehensive and novel hypotheses, addressing complex problems and accelerating the pace of discovery in today's evolving scientific landscape.

## 5 Categorization of Hypothesis Validation Approaches

The rapid advancements in hypothesis generation systems, particularly those powered by LLMs and computational intelligence, have highlighted the need for rigorous validation frameworks. Unlike traditional hypothesis formulation, where human intuition and prior knowledge play a central role, AI-generated hypotheses require systematic evaluation to ensure they are novel, insightful, scientifically plausible, testable, and actionable. Effective validation mechanisms serve as a critical checkpoint in the research lifecycle, distinguishing viable hypotheses from speculative or erroneous propositions.
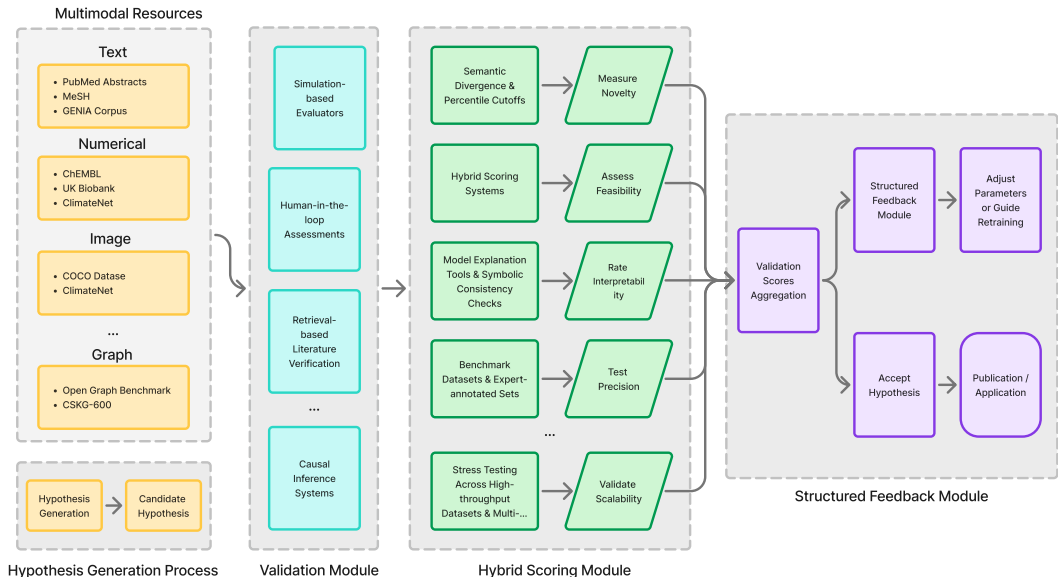


Figure 4: Pipeline for AI-assisted hypothesis validation. The figure outlines multiple validation modules, including simulation, human-in-the-loop assessments, retrieval-based verification, and causal inference, used to assess novelty, feasibility, interpretability, precision, and scalability. Finally, structured feedback and score aggregation are performed to guide acceptance, refinement, or retraining of hypotheses.

As scientific inquiries grow increasingly complex and interdisciplinary, traditional validation method-ologies, such as empirical testing and statistical inference, have been augmented by simulation-driven analysis, predictive modeling, causal inference techniques, and collaborative validation systems. These modern approaches integrate real-time data streams, machine learning models, and multi-agent validation frameworks, offering more scalable, adaptive, and domain-specific verification strategies. The convergence of computational methods and empirical validation represents a paradigm shift in how hypotheses are evaluated across diverse scientific fields, including biomedicine, social sciences, materials science, and artificial intelligence research. Figure 4 presents a multi-module architecture for hypothesis validation.

This section provides a structured taxonomy of hypothesis validation approaches, categorizing them based on their underlying methodologies, applications, and evaluation metrics. We explore how these approaches enhance reliability, reproducibility, and domain adaptability, ensuring that hypotheses meet the rigorous standards of contemporary scientific research. By examining these

validation frameworks, we highlight the evolving role of AI-assisted validation in fostering credible and impactful discoveries across disciplines.

Table 6: Summary of Hypothesis Validation Approaches, Tools, and Metrics

| Approach | Example Tool/System | Domain | Strengths | Metrics Used |
|---|---|---|---|---|
| Experimental and Empirical Testing | LabKey Server LabKey [2025] | Biomedical Research | Facilitates management and analysis of experimental data | Accuracy, Reproducibility Rate, Statistical Significance |
| Simulation-Driven Validation | Simulink MathWorks [2025] | Engineering, Robotics | Enables simulation of dynamic systems for hypothesis testing | Simulation Accuracy, Robustness Metric, Convergence Rate |
| Predictive, Adaptive, and Real-Time Validation | scikit-learn Pedregosa et al. [2011], Prophet Taylor and Letham [2018], TensorFlow Serving TensorFlow [2025] | Machine Learning, AI Research | Uses predictive models, real-time adaptation, and statistical inference for evolving datasets | Prediction Accuracy, Model Fit ($R^2$), Latency, Adaptation Accuracy |
| Cross-Disciplinary Generalization | IBM Watson IBM [2025] | Various Industries | Applies AI to ensure hypotheses are generalizable across domains | Transfer Learning Accuracy, Domain Alignment Score |
| Human-AI and Crowdsourced Validation | Zooniverse The Zooniverse Team [2025], Jupyter Notebook Project Jupyter [2025] | Citizen Science, AI, Education | Leverages human-AI interaction and crowdsourced expertise for validation | Agreement Rate, Consensus Score, User Satisfaction Metric, Iteration Success Rate |
| Causal Relationship Validation | TETRAD Center for Causal Discovery [2025] | Social Science, Medicine | Discovers causal relationships to validate hypotheses | Causal Effect Estimate, Goodness-of-Fit |
| Benchmarking and Standardized Testing | MLPerf MLPerf [2025] | Machine Learning | Provides benchmarks for evaluating model performance | F1-Score, Comparative Improvement |
| Multi-Agent Validation | SciAgents Ghafarollahi and Buehler [2024b] | AI, Computational Sciences | Leverages agent-based collaboration for scalable validation | Multi-Agent Consensus, Validation Confidence Score |
| Explainability and Interpretability Validation | SHAP Lundberg and Lee [2017], LIME Ribeiro et al. [2016] | AI, Social Sciences | Ensures interpretability of AI-generated hypotheses | Explainability Score, Feature Importance |

## 5.1 Experimental and Empirical Testing

Experimental validation remains a cornerstone of hypothesis testing, providing high reliability and precision through controlled experimentation. By manipulating variables and systematically observing their effects, this approach offers empirical evidence to confirm or refute hypotheses. It is particularly critical in fields such as biomedicine, chemistry, and physics, where direct causal inference is necessary to validate theoretical models Kim and Segev [2018]. Unlike computational approaches that rely on probabilistic modeling, experimental validation ensures direct empirical verification, making it indispensable for scientific rigor.

**Techniques and Tools.** Experimental validation employs several well-established methodologies to ensure accuracy, reproducibility, and generalizability. Controlled experiments minimize confounding factors by isolating specific variables, allowing precise measurement of their effects. Replication studies confirm the reliability of findings through repeated experimentation, ensuring consistency across different trials. Randomized Controlled Trials (RCTs) are a gold standard in biomedical research, enabling unbiased hypothesis testing by randomly assigning subjects to experimental conditions Beltagy et al. [2019].

Integrating automation and AI-driven experimental platforms has significantly enhanced the efficiency of traditional workflows. LabKey Server is widely adopted in biomedical research, facilitating data management and real-time analysis for experimental validation LabKey [2025]. In chemistry, robotic platforms integrated with AI-based automation improve reaction optimization and reduce human intervention, accelerating hypothesis testing Zhou et al. [2024]. Automated microscopy systems in experimental biology enable real-time cellular imaging and validation, minimizing manual errors while enhancing experimental throughput Qi et al. [2024].

**Metrics for Validation.** To ensure robust evaluation of experimental outcomes, multiple quantitative metrics are employed. Accuracy measures the correctness of experimental results by determining the proportion of correctly identified true positives and true negatives. Precision assesses the consistency of results across repeated experiments. The reproducibility rate ensures that independent researchers can validate findings, thereby strengthening scientific credibility Sybrandt et al. [2018]. Statistical significance, often determined through $p$-values ($p < 0.05$), verifies that observed effects are unlikely due to random chance, further supporting hypothesis validation Beltagy et al. [2019].

**Practical Applications.** Experimental validation is pivotal across scientific disciplines, ensuring that hypotheses translate into reliable real-world insights. In chemistry, AI-driven robotic laboratories optimize reaction conditions for synthetic compound development, improving reaction efficiency and hypothesis testing Zhou et al. [2024]. LabKey Server LabKey [2025] supports large-scale validation of biomarker discoveries in biomedicine, contributing to advancements in precision medicine. In physics, experimental validation through particle accelerators facilitates empirical testing of subatomic particle interactions, providing direct evidence for theoretical frameworks Kim and Segev [2018]. These applications highlight the essential role of experimental validation in bridging the gap between theoretical research and empirical confirmation.

**Feasibility and Novelty Accessibility.** Despite its strengths, experimental validation is often resource-intensive and time-consuming, requiring significant investment in specialized equipment, controlled environments, and expert personnel. These constraints can limit scalability, particularly for large-scale hypothesis testing or research in resource-constrained settings Qi et al. [2024]. Researchers increasingly integrate computational simulations and machine learning models with experimental validation to enhance feasibility, creating hybrid approaches that optimize resources while maintaining scientific rigor Zhou et al. [2024].

Regarding novelty accessibility, experimental validation is a key enabler of breakthrough discoveries, offering empirical verification of innovative hypotheses. Unlike predictive models that rely on inferential reasoning, experimental validation provides direct observational evidence, ensuring that novel hypotheses are rigorously tested before broader adoption. However, the slow throughput of traditional experimental techniques may hinder the rapid exploration of unconventional hypotheses. Integrating automated high-throughput screening platforms and AI-assisted experimental workflows helps overcome this limitation, allowing researchers to efficiently test and validate a larger number of hypotheses Sybrandt et al. [2018].

Experimental validation remains essential to scientific discovery, ensuring that hypotheses are empirically tested and reproducible. With the integration of automation, AI-driven validation, and hybrid computational-experimental approaches, this methodology continues to evolve, addressing the growing complexities of modern scientific inquiry while maintaining scientific rigor and reliability.

## 5.2 Simulation-Driven Validation

Simulation-based validation employs advanced computational models to emulate real-world systems, offering a scalable and risk-free environment for hypothesis testing. This approach is particularly valuable for hypotheses that involve high costs, extended timeframes, or significant safety risks if tested through physical experiments. By iteratively modeling and analyzing complex systems in silico, simulations enable the precise refinement and validation of hypotheses, ensuring their feasibility and reliability before practical implementation Kim and Segev [2018].

A key strength of simulation-based validation lies in its ability to mitigate physical risks and reduce the resources required for experimentation. Controlled virtual environments allow researchers to test hypotheses without physical constraints, eliminating the need for expensive and time-consuming laboratory setups. For example, automotive crash simulations significantly reduce costs and safety risks compared to physical crash tests. Moreover, simulations scale effectively to accommodate large

and multi-dimensional hypothesis spaces, enabling the simultaneous testing of numerous variables. This scalability and efficiency make simulation-driven validation an indispensable tool in modern scientific inquiry Beltagy et al. [2019].

**Techniques and Tools.** Simulation-based validation relies on several computational techniques to ensure robust hypothesis testing. Agent-Based Modeling (ABM) simulates the behavior of autonomous agents interacting in complex systems, making it highly applicable to ecology, economics, and social sciences Qi et al. [2024]. Monte Carlo Simulations use stochastic sampling to assess probabilistic outcomes, enabling the evaluation of hypothesis robustness under uncertainty Kim and Segev [2018]. Finite Element Analysis (FEA) is widely used in engineering and materials science to model structural and mechanical properties, validating physical hypotheses related to stress distribution, material behavior, and mechanical failure Sybrandt et al. [2018].

Several cutting-edge tools facilitate simulation-driven hypothesis validation. Simulink is a widely adopted platform for engineering and robotics, allowing researchers to model dynamic systems and test control strategies under simulated conditions MathWorks [2025]. BioSimulators, an advanced simulation toolkit for biomedical research, validates hypotheses related to molecular dynamics, enzyme kinetics, and cellular interactions LabKey [2025]. CARLA, an open-source autonomous vehicle simulator, provides a virtual environment for testing hypotheses related to robotic navigation, traffic dynamics, and urban mobility Zhou et al. [2024]. These tools exemplify the versatility and scalability of simulation-based validation, making it a powerful method for evaluating hypotheses across scientific disciplines.

**Metrics for Validation.** Simulation-based approaches employ several quantitative metrics to assess the validity, reliability, and robustness of simulated hypotheses. Simulation Accuracy measures how closely simulation outcomes align with empirical data, ensuring that the models effectively replicate real-world phenomena Beltagy et al. [2019]. Convergence Rate evaluates how quickly simulations stabilize, reflecting the efficiency of iterative model refinement. Robustness Metrics assess how hypotheses perform across diverse simulated conditions, ensuring reliability under different parameter settings Kim and Segev [2018]. These metrics provide a comprehensive framework for evaluating the precision and effectiveness of simulation-based validation.

**Practical Applications.** Simulation-driven validation is widely applied across scientific and engineering domains, offering researchers a safe and scalable platform for hypothesis testing. In biomedicine, tools like BioSimulators validate hypotheses about drug interactions, enzyme activity, and disease progression, allowing researchers to refine models before clinical trials LabKey [2025]. In robotics, CARLA enables the simulation of autonomous vehicle navigation, testing hypotheses related to sensor fusion, real-time decision-making, and safety protocols Zhou et al. [2024]. In materials science, Finite Element Analysis (FEA) is employed to validate stress-strain predictions, ensuring that hypotheses regarding mechanical durability and performance hold under real-world conditions Sybrandt et al. [2018]. These examples underscore the transformative role of simulation-based validation in hypothesis testing, providing risk-free and computationally scalable solutions.

**Feasibility and Novelty Accessibility.** Simulation-based validation is highly feasible, as it offers cost-effective and risk-free testing environments. Unlike laboratory-based experiments, simulations do not require physical resources, enabling researchers to iteratively refine models with minimal material costs Qi et al. [2024]. However, the effectiveness of this method is contingent upon the accuracy and completeness of the underlying models. Simplifications and assumptions in simulations may introduce biases, potentially limiting their ability to capture real-world complexities Kim and Segev [2018].

In terms of novelty accessibility, simulation-driven validation provides a safe and flexible environment for testing highly innovative and unconventional hypotheses. Unlike empirical testing, which requires physical prototypes and controlled conditions, simulations enable hypothesis exploration at scale, allowing researchers to test novel concepts that may not yet be feasible for real-world implementation Beltagy et al. [2019]. Nevertheless, the interpretability of simulation results remains a challenge, as complex computational models often require post-validation through empirical experiments to ensure external validity Sybrandt et al. [2018].

Simulation-driven validation continues to evolve as an indispensable methodology in modern scientific research. By integrating advanced computational techniques, AI-driven automation, and real-time predictive modeling, simulation-based validation offers a scalable, efficient, and innovative framework

for hypothesis testing. As computational models grow increasingly sophisticated, their role in scientific discovery and technological innovation will continue to expand, ensuring that simulation-driven validation remains a cornerstone of modern hypothesis evaluation.

## 5.3 Predictive, Adaptive, and Real-Time Validation

Predictive, adaptive, and real-time validation leverages machine learning models, statistical inference, and dynamic adaptation to evaluate hypotheses in a scalable, cost-effective, and data-driven manner. Unlike experimental approaches that require physical testing, these methods validate hypotheses through predictive modeling, real-time adjustments, and adaptive learning mechanisms. These approaches are particularly valuable for large-scale, complex, or continuously evolving datasets, making them indispensable in biomedicine, climate modeling, AI research, and financial forecasting Bengio et al. [2019].

A defining characteristic of predictive and adaptive validation is its ability to integrate high-dimensional datasets and continuous feedback loops, refining hypotheses and models iteratively. By enabling real-time hypothesis evaluation, these approaches dynamically adjust predictions as new data emerges, ensuring up-to-date and context-aware validation. This adaptability is crucial in fields where real-world conditions change rapidly, such as autonomous systems, financial markets, and environmental monitoring Beltagy et al. [2019].

**Techniques and Tools.** Predictive validation encompasses multiple modeling techniques to ensure robust hypothesis evaluation. Supervised learning models, widely applied in biomedicine and AI research, rely on labeled datasets for pattern recognition and classification Qi et al. [2024]. Bayesian inference models integrate prior knowledge with observed data, facilitating probabilistic hypothesis validation and dynamic uncertainty estimation, particularly in social sciences and healthcare analytics Kim and Segev [2018]. Time-series forecasting models, commonly used in climate science and financial analysis, predict trends and validate hypotheses based on historical and real-time data Zhou et al. [2024].

Several tools exemplify the application of these techniques. Scikit-learn provides robust validation frameworks for predictive modeling across various disciplines Pedregosa et al. [2011]. Prophet, a forecasting tool developed by Meta, specializes in time-series prediction and anomaly detection, aiding real-time hypothesis validation in economics and environmental sciences Taylor and Letham [2018]. TensorFlow Serving enables adaptive hypothesis validation in AI models, allowing real-time inference updates as new data becomes available TensorFlow [2025]. These tools highlight the versatility of predictive and adaptive validation, ensuring that hypotheses evolve dynamically with incoming data streams.

**Metrics for Validation.** To ensure rigorous evaluation, predictive validation relies on well-defined quantitative metrics. Prediction Accuracy assesses the correctness of model predictions by comparing forecasted outcomes with observed data Beltagy et al. [2019]. Bayesian Posterior Probability evaluates the likelihood of a hypothesis given available evidence, integrating both prior knowledge and new observations Kim and Segev [2018]. Model Fit ($R^2$) measures how well predictive models explain the variance in real-world datasets, serving as a key indicator of reliability Pedregosa et al. [2011]. Additionally, Adaptation Accuracy evaluates how effectively real-time models adjust to evolving data streams, ensuring continuous model refinement and reliability Bengio et al. [2019].

**Practical Applications.** Predictive, adaptive, and real-time validation plays a crucial role across scientific domains, enabling scalable and efficient hypothesis testing. In genomics, predictive models identify gene-disease associations, accelerating hypothesis validation in personalized medicine Qi et al. [2024]. In climate science, time-series forecasting models, such as Prophet, predict long-term temperature variations, aiding the validation of hypotheses about climate change and environmental feedback loops Taylor and Letham [2018]. In AI research, TensorFlow Serving supports adaptive model validation, ensuring that machine learning models remain accurate as they process real-time streaming data TensorFlow [2025]. These applications underscore the significance of predictive validation in handling complex, dynamic, and continuously evolving scientific challenges.

**Feasibility and Novelty Accessibility.** Predictive validation offers high feasibility due to its cost-effectiveness and scalability. Unlike traditional experimental validation, which requires physical setups and manual intervention, predictive approaches allow researchers to test and refine hypotheses iteratively using existing datasets and real-time observations. However, the effectiveness of predictive

validation depends on the availability of high-quality training data and the interpretability of complex models Zhou et al. [2024]. Poorly curated datasets or biased training samples can lead to overfitting, reducing the generalizability of validated hypotheses Beltagy et al. [2019].

In terms of novelty accessibility, predictive validation enables hypothesis exploration in high-dimensional spaces, revealing patterns and relationships that might be overlooked in traditional experimental settings. This approach is particularly valuable in AI-assisted discovery, where deep learning models generate hypotheses that human researchers might not have considered Sybrandt et al. [2018]. However, the insights generated by predictive models are often constrained by the assumptions inherent in the models themselves, limiting their ability to fully explore unconventional or outlier hypotheses. As a result, hybrid validation approaches—combining predictive modeling with empirical validation—are increasingly being adopted to enhance both scalability and interpretability.

Predictive, adaptive, and real-time validation continues to evolve as an essential methodology for hypothesis testing, bridging the gap between static theoretical models and dynamic real-world applications. As advancements in machine learning, real-time analytics, and automated hypothesis refinement progress, predictive validation will play an increasingly critical role in accelerating scientific discovery and technological innovation.

## 5.4   Cross-Disciplinary Generalization

Cross-domain validation is a powerful approach for assessing hypotheses by evaluating their applicability across multiple scientific fields or datasets. This methodology ensures that hypotheses are generalizable and robust, making it particularly valuable in interdisciplinary research where knowledge transfer and cross-disciplinary correlations are essential. By testing hypotheses in diverse domains, cross-domain validation facilitates the discovery of broader applications and fosters the integration of knowledge across scientific disciplines.

A key strength of cross-domain validation lies in its ability to encourage generalization by testing hypotheses across varying fields. This approach not only ensures the relevance of hypotheses beyond their original context but also enables the discovery of novel insights that span multiple domains. Additionally, cross-domain validation facilitates the transfer of knowledge and methodologies, providing a framework for integrating diverse datasets and experimental techniques, thereby strengthening the coherence and applicability of hypotheses.

**Techniques and Tools.** The effectiveness of cross-domain validation is supported by several innovative techniques. *Domain Mapping* aligns datasets and methodologies across fields to ensure consistency and compatibility during validation, addressing challenges of heterogeneity in cross-domain data Zhang et al. [2024b]. *Interdisciplinary Networks* utilize network models to identify shared structures and relationships, enhancing the coherence of hypotheses by uncovering commonalities across scientific disciplines Sybrandt et al. [2018]. *Transfer Learning* is another critical technique, applying knowledge derived from one domain to test and validate hypotheses in another, enabling adaptability and scalability in hypothesis evaluation Wu et al. [2020].

Several tools and systems exemplify the application of these techniques. *CrossValNet* integrates datasets from genomics, pharmacology, and environmental science, enabling researchers to validate hypotheses that span these fields Sybrandt et al. [2018]. *InterdisciplinaryTest* facilitates hypothesis validation across domains such as physics, chemistry, and biology, ensuring compatibility and coherence in results Zhou et al. [2024]. *TransferTest* employs transfer learning techniques to adapt and validate AI models and hypotheses across diverse scientific domains, highlighting its versatility and impact Touvron et al. [2023].

**Metrics for Validation.** Cross-domain validation employs several metrics to assess the robustness and applicability of hypotheses. *Domain Alignment Score* measures the consistency of results across different domains, providing a quantitative evaluation of hypothesis generalizability. *Transfer Learning Accuracy* evaluates the effectiveness of knowledge transfer during validation, ensuring that models and hypotheses perform reliably in target domains. *Interdisciplinary Correlation Metric* assesses the strength of relationships between variables in different fields, enabling researchers to identify and quantify cross-disciplinary connections. These metrics collectively provide a comprehensive framework for evaluating the success of cross-domain validation efforts.

**Practical Applications.** The practical applications of cross-domain validation span a wide range of scientific fields, underscoring its adaptability and utility. In genomics and pharmacology, *CrossVal-Net* has validated gene-disease associations by integrating genomic and pharmacological datasets, demonstrating the broader relevance of these hypotheses Sybrandt et al. [2018]. In environmental science and biology, *InterdisciplinaryTest* has explored connections between ecological variables and biological diversity, providing actionable insights into conservation strategies Zhou et al. [2024]. In AI and robotics, *TransferTest* has validated hypotheses about navigation algorithms by transferring models from robotics to autonomous vehicles, showcasing its ability to ensure adaptability across diverse scientific fields Touvron et al. [2023].

**Feasibility and Novelty Accessibility.** Cross-domain validation is highly feasible for fostering hypothesis robustness and facilitating the integration of diverse datasets and methodologies. Its strengths lie in its capacity to generalize hypotheses across multiple domains, ensuring their reliability and applicability. However, the approach requires extensive expertise across different fields and significant efforts to align heterogeneous datasets, which can pose challenges.

In terms of novelty, cross-domain validation excels in facilitating the discovery of novel connections and insights by leveraging domain-specific knowledge. By integrating datasets and methodologies from diverse disciplines, this approach uncovers relationships that might otherwise remain hidden. However, the availability and compatibility of cross-domain datasets may limit its broader applications, particularly in emerging or underexplored fields.

Cross-domain validation represents a transformative methodology for hypothesis testing in interdisciplinary research. By encouraging generalization, uncovering cross-disciplinary correlations, and fostering knowledge transfer, this approach continues to drive innovation and discovery across diverse scientific landscapes.

## 5.5   Human-AI and Crowdsourced Validation

Human-AI and crowdsourced validation leverage both computational efficiency and collective intelligence to test and evaluate hypotheses. This combined approach democratizes hypothesis validation by integrating human intuition, domain expertise, and machine learning capabilities, enabling large-scale, iterative, and context-aware evaluations. Crowdsourced validation engages diverse participant groups to provide broad insights, while Human-AI collaboration refines hypotheses through interactive feedback loops and explainable AI, enhancing interpretability and robustness. These methodologies are particularly beneficial in domains requiring broad engagement, adaptability, and transparency.

A defining strength of this approach is its ability to combine the scalability of crowdsourced validation with the precision and adaptability of Human-AI collaboration. Crowdsourced validation ensures that hypotheses are tested across a wide range of perspectives, while AI-driven refinements optimize accuracy and efficiency. Furthermore, iterative feedback loops foster continuous hypothesis improvement, ensuring that both human insights and machine learning capabilities contribute dynamically to the validation process.

**Techniques and Tools.** The effectiveness of Human-AI and crowdsourced validation is supported by several innovative techniques. *Open Feedback Systems* collect real-time user feedback, enabling dynamic adjustments during hypothesis validation Zhou et al. [2024]. *Gamification* engages participants through interactive elements to sustain engagement and improve data quality Kim and Segev [2018]. *Consensus Aggregation* synthesizes diverse participant inputs to derive statistically robust conclusions, mitigating outlier influence and ensuring reliability Sybrandt et al. [2018]. Additionally, *Interactive AI Feedback* facilitates iterative refinements by allowing users to guide AI-driven hypothesis improvements, ensuring adaptability and contextual relevance Zhou et al. [2024].

Several tools exemplify the application of these techniques. *Zooniverse* provides a crowdsourced platform where participants contribute to scientific discovery by validating hypotheses across various disciplines The Zooniverse Team [2025]. *Jupyter Notebook* facilitates collaborative hypothesis testing by integrating human expertise with AI-driven analytics Project Jupyter [2025]. *Explain-Valid* combines explainable AI techniques with user feedback, improving transparency and trust in hypothesis validation processes Zhou et al. [2024]. *FeedbackLoopAI* integrates iterative user feedback into machine learning models, ensuring that hypotheses evolve dynamically based on human insights Sybrandt et al. [2018].

**Metrics for Validation.** To ensure the reliability and effectiveness of Human-AI and crowdsourced validation, multiple quantitative metrics are employed. The *Agreement Rate* measures response consistency across participants, providing insight into hypothesis reliability. The *Consensus Score* quantifies the degree of collective agreement, ensuring robust validation through diverse perspectives. The *User Engagement Metric* evaluates the level of participant involvement, indicating the scalability and accessibility of the validation process. In Human-AI interactions, the *Iteration Success Rate* assesses how effectively hypotheses are refined over successive validation rounds, while the *Explainability Score* quantifies the interpretability of AI-generated outputs based on human verification Zhou et al. [2024].

**Practical Applications.** The integration of Human-AI and crowdsourced validation has demonstrated significant utility across scientific fields. In social sciences, *Zooniverse* has validated hypotheses about public perception and social behaviors through large-scale participant contributions The Zooniverse Team [2025]. In Human-Computer Interaction (HCI), *Jupyter Notebook* has facilitated usability testing by integrating human insights with machine learning-driven analytics Project Jupyter [2025]. In AI ethics, *ExplainValid* has enabled researchers to validate ethical guidelines by offering transparent, interpretable AI-generated recommendations Zhou et al. [2024]. Additionally, in healthcare, *FeedbackLoopAI* has assisted clinicians in refining diagnostic models through iterative AI-assisted hypothesis validation Sybrandt et al. [2018].

**Feasibility and Novelty Accessibility.** Human-AI and crowdsourced validation offer high feasibility due to their scalability, cost-effectiveness, and accessibility. Crowdsourced validation provides diverse perspectives while minimizing resource constraints, making it particularly useful for hypothesis testing in large and distributed participant groups. However, challenges such as demographic biases and response variability must be managed to ensure generalizability. In contrast, Human-AI collaboration enhances feasibility in complex, domain-specific tasks requiring expert intervention, though it may be limited by scalability constraints.

In terms of novelty, this hybrid approach excels at uncovering unconventional hypotheses through broad participant contributions and AI-driven insights. Crowdsourced validation fosters creativity by integrating diverse viewpoints, while Human-AI collaboration enhances the precision and interpretability of results. However, reliance on human input can introduce biases, and AI models may struggle with outlier hypotheses that deviate from learned patterns. Addressing these limitations requires careful design of validation frameworks that optimize both human intuition and machine-driven analysis.

Human-AI and crowdsourced validation represent a transformative approach to hypothesis testing, combining the power of collective intelligence with computational efficiency. By leveraging iterative feedback, interdisciplinary collaboration, and scalable participation, this approach fosters inclusivity, adaptability, and rigor across a wide range of scientific and practical applications.

## 5.6 Causal Relationship Validation

Causal relationship validation is a critical methodology for establishing cause-and-effect relationships between variables, distinguishing genuine causal mechanisms from mere correlations. This approach is particularly valuable in fields such as social sciences, biomedicine, and economics, where controlled experiments may be impractical or infeasible. By leveraging both observational and experimental data, causal inference techniques provide a structured framework for hypothesis validation, allowing for deeper insights into underlying mechanisms.

A key strength of causal relationship validation lies in its focus on causality rather than associative patterns, ensuring that validated hypotheses have mechanistic explanations beyond statistical relationships. This methodological rigor enhances the applicability of hypotheses in real-world decision-making processes, making it essential in fields requiring strong evidence-based conclusions. Additionally, causal inference techniques bridge the gap between observational studies and actionable insights, reinforcing their importance in modern scientific research.

**Techniques and Tools.** Causal inference validation employs several advanced methodologies to establish reliable causal relationships. *Structural Causal Models (SCMs)* use directed acyclic graphs (DAGs) to explicitly represent causal pathways, enabling systematic validation of mechanistic hypotheses Pearl [2009]. *Propensity Score Matching (PSM)* reduces confounding effects by simulating randomization in observational datasets, improving the robustness of causal conclusions Rosenbaum

and Rubin [1983]. *Instrumental Variable Analysis (IVA)* introduces external instruments to estimate causal effects when direct experimentation is infeasible, making it particularly useful in econometric studies Angrist et al. [1996].

Several computational tools facilitate the practical implementation of these techniques. *TETRAD*, a widely used causal discovery platform, applies SCMs to infer causal structures in domains such as epidemiology and behavioral sciences Center for Causal Discovery [2025]. *Scikit-learn*, though primarily a machine learning library, includes causal inference modules that support propensity score-based methods for causal validation Pedregosa et al. [2011]. *IBM Watson* employs AI-driven causal modeling techniques to validate hypotheses across multiple industries, offering insights into customer behavior, healthcare diagnostics, and economic forecasting IBM [2025].

**Metrics for Validation.** Causal inference validation employs several quantitative metrics to assess the strength and reliability of causal relationships. The *Causal Effect Estimate* quantifies the magnitude of causal influences by comparing outcomes across treated and control groups, ensuring precise hypothesis evaluation Rosenbaum and Rubin [1983]. The *Goodness-of-Fit for SCMs* evaluates how well causal models explain observed data, offering a quantitative measure of model reliability Pearl [2009]. Additionally, the *Instrument Relevance Score* assesses the validity of instrumental variables, ensuring that they provide unbiased causal estimates and meet key econometric assumptions Angrist et al. [1996].

**Practical Applications.** Causal inference validation has been successfully applied across a range of disciplines. In biomedical research, *TETRAD* has been instrumental in identifying causal links between genetic variants and disease phenotypes, advancing precision medicine Center for Causal Discovery [2025]. In policy evaluation, *IBM Watson* has been used to assess the causal impact of education and healthcare policies, aiding policymakers in evidence-based decision-making IBM [2025]. In econometrics, *Scikit-learn*'s causal inference techniques have facilitated the evaluation of fiscal policies and labor market trends, supporting data-driven economic strategies Pedregosa et al. [2011].

**Feasibility and Novelty Accessibility.** Causal relationship validation offers high feasibility due to its adaptability to both observational and experimental data. Its ability to infer causality from complex datasets makes it a valuable tool across multiple domains. However, challenges such as confounding bias, instrument validity, and model specification errors must be carefully managed to ensure accurate results. The effectiveness of causal inference methods depends on the availability of high-quality data and the correct selection of causal assumptions, which can be limiting factors in some studies.

In terms of novelty, causal inference validation excels at uncovering previously unidentified causal pathways, enriching scientific understanding. By providing mechanistic explanations, it extends hypothesis validation beyond surface-level correlations, enabling deeper insights into fundamental processes. However, causal inference models rely on assumptions such as the validity of instrumental variables and the correctness of DAG structures, which may constrain the exploration of novel hypotheses. Advances in hybrid causal discovery techniques and the integration of AI-driven causal modeling promise to further enhance the scope and accessibility of causal inference validation.

Causal relationship validation remains a cornerstone of modern scientific research, offering robust methodologies to uncover the fundamental drivers of complex phenomena. As computational techniques and data integration frameworks continue to evolve, causal inference will play an increasingly critical role in refining hypotheses and guiding evidence-based decision-making across diverse scientific disciplines.

## 5.7 Benchmarking and Standardized Testing

Benchmarking and standardized testing serve as foundational approaches for hypothesis evaluation, utilizing predefined datasets, protocols, and metrics to ensure reliability, accuracy, and comparability. This methodology is particularly prevalent in fields such as machine learning, artificial intelligence, and experimental sciences, where the ability to systematically compare hypotheses against established baselines fosters reproducibility and facilitates objective assessment. By providing quantifiable measures of performance, benchmarking enables researchers to identify strengths, weaknesses, and areas for improvement in their hypotheses, accelerating scientific and technological advancements.

A key characteristic of benchmarking is its emphasis on reproducibility and consistency, ensuring that hypotheses can be evaluated across diverse research environments. Standardized testing offers a structured framework for assessing performance relative to existing solutions, enabling transparent and evidence-based decision-making. This comparative approach also helps identify the competitive advantages of new hypotheses, facilitating broader adoption in real-world applications.

**Techniques and Tools.** Benchmarking and standardized testing employ a range of established techniques to ensure rigorous validation. *Baseline Comparison* assesses hypothesis performance against predefined reference models, providing a direct measure of improvement MLPerf [2025]. *Cross-Validation*, widely used in machine learning and experimental sciences, partitions data into training and testing subsets to ensure robustness and generalizability Pedregosa et al. [2011]. *Performance Metrics Analysis* applies standardized quantitative measures, such as accuracy, F1-score, and comparative improvement, to comprehensively evaluate hypotheses MLPerf [2025].

Several tools exemplify these techniques in various domains. *MLPerf*, an industry-standard benchmarking suite, provides reproducible evaluations of machine learning models across multiple tasks, enabling rigorous hypothesis testing MLPerf [2025]. *Scikit-learn*, a widely used machine learning library, includes built-in benchmarking utilities for evaluating model generalization across datasets Pedregosa et al. [2011]. *Simulink*, a simulation platform, facilitates the comparative testing of control systems and engineering hypotheses under standardized experimental conditions MathWorks [2025].

**Metrics for Validation.** Standardized testing employs well-defined metrics to assess hypothesis performance objectively. *Accuracy* measures correctness by comparing predicted outcomes with observed results:

$$\text{Accuracy} = \frac{\text{Correct Predictions}}{\text{Total Predictions}}$$

The *F1-Score*, balancing precision and recall, is particularly useful for hypotheses evaluated on imbalanced datasets:

$$F1 = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

Additionally, *Comparative Improvement* quantifies the relative performance gain of a hypothesis over baseline solutions:

$$\text{Improvement} = \frac{\text{Hypothesis Performance} - \text{Baseline Performance}}{\text{Baseline Performance}}$$

**Practical Applications.** Benchmarking and standardized testing have demonstrated their utility across a broad range of disciplines. In machine learning, *MLPerf* has validated neural network architectures by benchmarking their performance on datasets such as ImageNet, driving advancements in model optimization and evaluation MLPerf [2025]. In biomedical science, *Scikit-learn* has supported hypothesis validation in genomics and disease prediction, ensuring reproducibility in clinical applications Pedregosa et al. [2011]. In engineering and robotics, *Simulink* has benchmarked control algorithms and navigation models under standardized experimental conditions, facilitating the development of reliable autonomous systems MathWorks [2025].

**Feasibility and Novelty Accessibility.** The feasibility of benchmarking lies in its structured and reproducible validation framework, which accelerates the adoption of validated hypotheses by quantifying their advantages over existing solutions. However, its effectiveness depends on the availability and relevance of benchmark datasets, which may become outdated or fail to capture emerging trends, limiting their applicability.

In terms of novelty, benchmarking effectively highlights innovations by measuring hypotheses relative to established baselines. This ensures that advancements are grounded in objective comparisons and reproducible evidence. However, reliance on predefined benchmarks can sometimes constrain exploratory research, as benchmarks may favor incremental improvements over paradigm-shifting hypotheses.

Benchmarking and standardized testing provide a rigorous and structured methodology for hypothesis evaluation, ensuring reliability, comparability, and scientific rigor. By leveraging standardized datasets and performance metrics, this approach fosters reproducibility and drives progress across machine learning, experimental sciences, and engineering disciplines.

## 5.8 Multi-Agent Validation

Multi-agent validation leverages autonomous and collaborative agents to systematically test and refine hypotheses, providing a scalable and distributed approach to validation. By utilizing multiple intelligent agents that interact, exchange information, and adapt their strategies, this methodology enables efficient validation in complex, dynamic, and data-rich environments. Multi-agent validation is particularly valuable in artificial intelligence, computational sciences, and decision-making systems, where multiple perspectives and decentralized reasoning enhance the reliability and robustness of validated hypotheses.

A defining characteristic of multi-agent validation is its ability to simulate complex real-world interactions by distributing hypothesis evaluation across multiple agents. These agents operate in parallel, facilitating large-scale hypothesis validation while dynamically adjusting their strategies based on continuous feedback. This approach enhances robustness by reducing bias and increasing the diversity of perspectives incorporated into the validation process.

**Techniques and Tools.** Multi-agent validation employs a range of techniques to coordinate and optimize the validation process. *Consensus Mechanisms* ensure that agents reach agreement on hypothesis validity through iterative feedback and refinement, enhancing reliability Ghafarollahi and Buehler [2024b]. *Distributed Learning* enables agents to learn and adapt to new data independently while collectively improving hypothesis validation outcomes Ghafarollahi and Buehler [2024b]. *Reinforcement Learning-Based Coordination* optimizes agent collaboration, allowing them to efficiently explore and validate hypotheses in complex environments Ghafarollahi and Buehler [2024b].

Several computational frameworks support multi-agent validation across different domains. *SciAgents*, a multi-agent system designed for hypothesis generation and validation, coordinates autonomous agents to evaluate scientific claims with high efficiency Ghafarollahi and Buehler [2024b]. *Protagents* integrates multiple AI-driven agents to validate hypotheses in dynamic knowledge environments, improving adaptability and inference capabilities Ghafarollahi and Buehler [2024a]. *Agent-based Decision Platforms* utilize decentralized agents for hypothesis testing in strategic decision-making and simulation-driven research, allowing for scalable and real-time validation Ghafarollahi and Buehler [2024b].

**Metrics for Validation.** Multi-agent validation relies on specialized metrics to assess hypothesis consistency and agent performance. The *Multi-Agent Consensus Score* quantifies the degree of agreement among agents, ensuring reliability in hypothesis evaluation:

$$\text{Consensus Score} = \frac{\text{Number of Agents in Agreement}}{\text{Total Number of Agents}}$$

The *Validation Confidence Score* evaluates the confidence level of hypothesis acceptance, incorporating uncertainty estimates from multiple agents:

$$\text{Validation Confidence} = \frac{\sum \text{Agent Confidence Scores}}{\text{Number of Agents}}$$

Additionally, the *Exploration-Exploitation Balance* measures how well agents optimize hypothesis validation by balancing novel exploration with confirmatory validation:

$$\text{Exploration-Exploitation Ratio} = \frac{\text{Exploratory Actions}}{\text{Confirmatory Actions}}$$

**Practical Applications.** Multi-agent validation has demonstrated its effectiveness across diverse scientific and technological domains. In AI research, *SciAgents* has been used to validate machine learning-generated hypotheses by leveraging collaborative AI agents, ensuring improved generalization and robustness Ghafarollahi and Buehler [2024b]. In scientific discovery, *Protagents* has facilitated hypothesis validation by enabling autonomous knowledge exploration in interdisciplinary research Ghafarollahi and Buehler [2024a]. In computational sciences, multi-agent decision platforms have been applied to strategic planning and simulation-based hypothesis validation, improving decision accuracy in complex systems Ghafarollahi and Buehler [2024b].

**Feasibility and Novelty Accessibility.** Multi-agent validation offers high feasibility due to its ability to scale hypothesis validation processes efficiently. By distributing computational tasks among agents, this approach enables large-scale validation without requiring centralized control. However,

challenges such as agent coordination, communication overhead, and decision-making biases must be carefully managed to maintain validation accuracy.

In terms of novelty, multi-agent validation fosters innovative hypothesis exploration by leveraging independent but interconnected agents. This decentralized approach allows for the discovery of unconventional insights that may be overlooked in traditional validation frameworks. However, the reliance on multi-agent interactions introduces complexity, requiring robust mechanisms to manage conflicts, align validation strategies, and ensure interpretability.

Multi-agent validation continues to revolutionize hypothesis testing by integrating autonomous decision-making, distributed learning, and collaborative inference. As AI and computational sciences advance, multi-agent systems will play an increasingly critical role in ensuring scalable, efficient, and high-confidence hypothesis validation across scientific disciplines.

## 5.9 Explainability and Interpretability Validation

Explainability and interpretability validation ensures that hypotheses and their underlying models are not only accurate but also transparent and understandable. This approach is particularly crucial in domains such as artificial intelligence, social sciences, and biomedical research, where the ability to explain results fosters trust, reproducibility, and informed decision-making. By validating hypotheses through interpretable frameworks, this methodology enhances transparency, accountability, and the reliability of complex models.

A defining characteristic of explainability and interpretability validation is its focus on making complex hypothesis evaluation processes understandable to both domain experts and non-specialists. This is particularly important in AI-driven hypothesis generation, where black-box models can produce highly accurate but opaque results. Explainability validation helps bridge this gap by ensuring that the reasoning behind hypothesis acceptance or rejection is accessible and interpretable.

**Techniques and Tools.** Explainability and interpretability validation employ various techniques to clarify how hypotheses are derived and assessed. *Feature Attribution Methods* identify the most influential variables in hypothesis evaluation, ensuring transparency in model-driven hypothesis testing Lundberg and Lee [2017]. *Local Interpretable Model-Agnostic Explanations (LIME)* approximates black-box models with simpler, interpretable models to provide local explanations for hypothesis validation Ribeiro et al. [2016]. *Counterfactual Reasoning* generates "what-if" scenarios to test how slight modifications in input variables affect the hypothesis outcome, ensuring robustness and fairness Ribeiro et al. [2016].

Several tools facilitate explainability and interpretability validation. *SHAP (Shapley Additive Explanations)* provides game-theoretic explanations for complex models, quantifying feature importance in hypothesis validation Lundberg and Lee [2017]. *LIME* generates local surrogate models to make AI-generated hypotheses more interpretable Ribeiro et al. [2016]. *IBM Watson Explainability Tools* integrate explainability methods into AI-driven hypothesis testing, ensuring that predictions align with human reasoning and regulatory requirements IBM [2025].

**Metrics for Validation.** Explainability validation relies on well-defined metrics to assess interpretability and transparency. The *Explainability Score* quantifies the degree to which a hypothesis or model outcome can be understood by human users:

$$\text{Explainability Score} = \frac{\text{Number of Correctly Interpreted Predictions}}{\text{Total Predictions}}$$

The *Feature Importance Consistency* metric evaluates how consistently a model ranks the importance of different variables across different runs:

$$\text{Feature Consistency} = \frac{\text{Stable Feature Rankings across Runs}}{\text{Total Features Ranked}}$$

Additionally, the *Human Trust Index* measures user confidence in hypothesis validation results based on their clarity and interpretability:

$$\text{Trust Index} = \frac{\text{Number of Users Who Trust Explanations}}{\text{Total Users Surveyed}}$$

**Practical Applications.** Explainability and interpretability validation have demonstrated significant impact across diverse domains. In AI research, *SHAP* has been applied to validate machine learning

hypotheses by explaining the influence of different features on model predictions, improving model trustworthiness Lundberg and Lee [2017]. In biomedical science, *LIME* has been used to interpret AI-driven disease diagnosis hypotheses, ensuring that medical predictions align with clinical reasoning Ribeiro et al. [2016]. In regulatory compliance, *IBM Watson Explainability Tools* have facilitated transparent hypothesis validation in financial and legal decision-making, ensuring adherence to explainability mandates IBM [2025].

**Feasibility and Novelty Accessibility.** Explainability validation is increasingly feasible due to its widespread adoption in fields such as artificial intelligence, healthcare, and policymaking, where it enhances trust, transparency, and accessibility in hypothesis validation. By making model behavior more interpretable, explainability techniques support a clearer understanding and communication of computational reasoning. They also promote scientific innovation by revealing hidden patterns within hypotheses that can lead to deeper insights. However, challenges persist in maintaining a balance between interpretability and model performance, as overly simplified explanations may obscure the complexity of advanced hypotheses. Moreover, existing explainability tools often struggle with abstract or unconventional hypotheses generated by models that lack well-defined decision boundaries, limiting their applicability in frontier scientific contexts. Despite these limitations, explainability and interpretability remain foundational for promoting scientific rigor, especially as hypothesis generation becomes more data-driven and automated. Ensuring transparent and interpretable validation will be essential for bridging the gap between complex computational reasoning and actionable, real-world decision-making.

## 5.10 Hybrid Validation Methods

While individual validation strategies such as empirical testing, simulation, or predictive modeling are valuable independently, many scientific domains benefit most from hybrid validation methods that combine multiple approaches to enhance reliability, generalizability, and interpretability Qi et al. [2024], Aubin Le Quéré et al. [2024]. Hybrid validation frameworks leverage the complementary strengths of diverse validation techniques to mitigate individual limitations and address the complex nature of scientific hypotheses. For example, in biomedical research, hypotheses often undergo predictive validation using trained classifiers or generative models, followed by simulation-based testing in virtual biological environments Laurent et al. [2024]. Promising hypotheses are then escalated to experimental validation under laboratory conditions Schmidgall et al. [2024]. This progressive validation pipeline balances scalability with scientific rigor and helps conserve resources by filtering out implausible hypotheses early in the process. In materials science, hybrid frameworks often integrate simulation-driven methods with causal inference models to understand the mechanistic underpinnings of observed behaviors Ghafarollahi and Buehler [2024b]. High throughput simulations generate candidate hypotheses, while statistical methods validate causal relations among physical properties. These insights are subsequently verified through real-world experiments, enabling both theoretical insight and empirical confidence. Furthermore, human AI collaborative validation is frequently combined with benchmarking and explainability techniques. For instance, hypotheses generated by language models can be benchmarked against curated datasets or domain-specific baselines, and their justifications can be interpreted using explainable artificial intelligence methods Shavit et al. [2023]. Domain experts review these explanations to ensure alignment with established field standards. This hybrid setup ensures that hypotheses are computationally viable, epistemically sound, and ethically aligned Jaradeh et al. [2019].

Hybrid validation methods also support interdisciplinary hypothesis evaluation, where no single domain has sufficient ground truth or validation infrastructure. In such settings, modular validation architectures activate specialized components for different subdomains, such as simulations for physics, knowledge graphs for biology, and crowd-based assessment for sociotechnical hypotheses while maintaining a coherent end-to-end evaluation pipeline Zhou et al. [2024], Sybrandt et al. [2017]. Finally, hybrid validation facilitates risk-aware hypothesis evaluation. High-risk, high-impact ideas can be evaluated using risk-weighted scoring functions that combine quantitative metrics such as predictive accuracy and novelty with qualitative assessments such as feasibility narratives obtained from expert feedback Fok and Weld [2024], Wang et al. [2023]. These approaches are especially critical in frontier scientific research, where uncertainty is high and exploratory validation is necessary.

In summary, hybrid validation methods offer a pragmatic and flexible approach to hypothesis evaluation by synthesizing diverse validation signals into a coherent judgment framework. As scientific
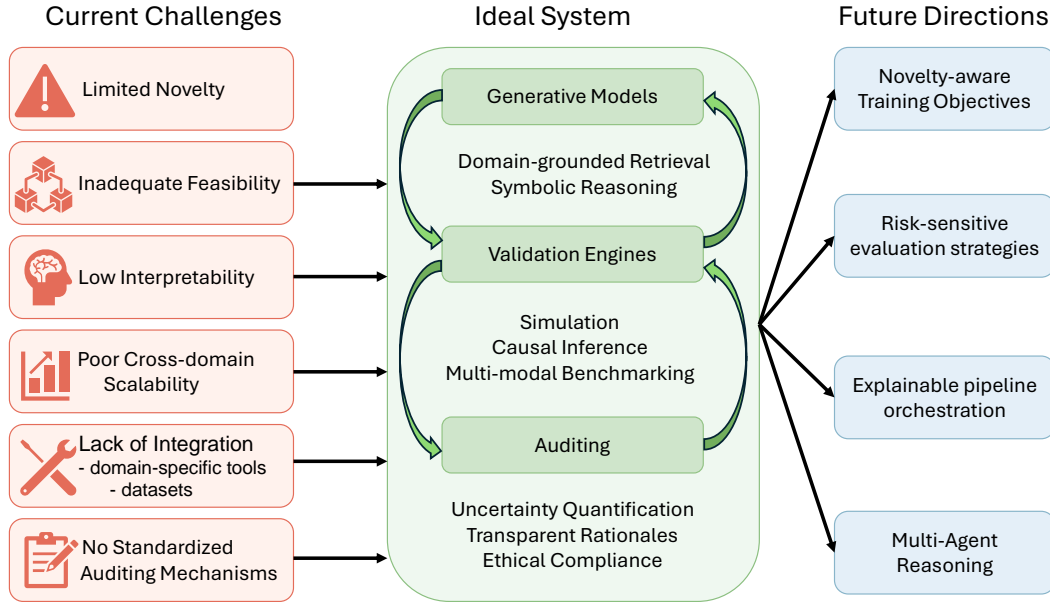
Figure 5: Roadmap for future directions in LLM-based scientific hypothesis generation and validation. The figure connects current challenges, such as limited novelty, feasibility issues, and lack of interpretability, to future directions like novelty-aware training objectives, risk-sensitive evaluation strategies, explainable pipeline orchestration, and multi-agent reasoning. The figure outlines components of an ideal system integrating generative models, validation engines, and auditing.

discovery becomes increasingly interdisciplinary and automated, hybrid validation pipelines will be essential for balancing scalability, trust, and scientific robustness.

**Conclusion.** The categorization of hypothesis validation methodologies highlights their diverse applicability across scientific and technological domains. Each approach, from the *precision and empirical rigor of experimental validation* to the *scalability and efficiency of predictive and simulation-driven models*, plays a distinct role in ensuring the robustness of validated hypotheses. *Benchmarking and standardized testing* provide structured comparability, fostering reproducibility and competitive assessment, while *causal inference validation* uncovers mechanistic relationships beyond correlations, strengthening the explanatory power of scientific claims. The integration of *human-AI collaboration and crowdsourced validation* enhances interpretability and accessibility, ensuring that hypotheses are not only statistically sound but also contextually meaningful. *Cross-disciplinary and multi-agent validation* extend these capabilities by enabling adaptability and scalability, allowing knowledge transfer across diverse domains and leveraging distributed intelligence for large-scale hypothesis testing. Furthermore, *explainability-focused validation approaches* enhance transparency and trust, bridging the gap between complex computational models and human decision-making.

By strategically combining these methodologies, researchers can construct *robust, scalable, and adaptive validation frameworks* tailored to modern, data-driven research challenges. This convergence of methodologies not only accelerates scientific discovery but also fosters *innovation, interdisciplinary breakthroughs, and higher standards of reliability and generalizability in hypothesis validation*, ensuring that research remains both impactful and reproducible in an ever-evolving scientific landscape.

# 6 Future Directions: Overcoming Limitations in Hypothesis Creation and Validation

The transformative potential of LLMs in scientific hypothesis generation and validation is undeniable. However, their current capabilities are constrained by challenges such as limited novelty generation, feasibility assessment, pervasive data biases, inadequate interpretability, scalability, and adaptability. Addressing these limitations is essential to fully harness their potential across diverse scientific domains. While LLM-driven advancements in hypothesis generation and validation have

shown significant promise, the complexity and interdisciplinary demands of modern research require solutions that transcend existing limitations. Realizing the potential of LLMs to address high-risk, cross-disciplinary challenges demands a focused effort to overcome technical and methodological barriers, ensuring their integration into diverse scientific workflows. Figure 5 provides a future roadmap connecting current limitations in LLM-driven hypothesis generation and validation with targeted capabilities needed in next-generation systems to achieve an ideal system.

This section highlights key challenges faced by LLM-based systems and presents actionable strategies to address them. By exploring emerging methodologies and identifying opportunities for innovation, we aim to provide a strategic roadmap for advancing hypothesis-driven discovery. These efforts will help make LLMs more robust, adaptable, and impactful tools for fostering scientific breakthroughs and addressing pressing global challenges.

## 6.1 Enhancing Novelty in Hypothesis Creation

Large language models (LLMs) are often constrained by their reliance on established knowledge bases and historical data, which can limit their ability to produce transformative or genuinely novel hypotheses. While this grounding ensures that hypotheses align with current scientific understanding, it also risks reinforcing conventional paradigms. Such reliance on past knowledge restricts the exploration of counterintuitive or uncharted areas, particularly in fields like theoretical sciences or frontier technologies that demand innovative thinking. Such reliance inadvertently reinforces existing biases and hinders transformative innovation. These limitations have several implications. They restrict the potential for paradigm-shifting discoveries, often favoring incremental advancements over groundbreaking ideas. Furthermore, fields that require unconventional or speculative approaches remain underserved, and underrepresented or unconventional research areas risk being overlooked, perpetuating existing biases in scientific exploration.

**Strategies for Overcoming Limitations:**

1. **Incorporating Generative Exploration Models:** Utilize creative frameworks, such as generative adversarial networks (GANs) or reinforcement learning, to introduce variability and randomness into hypothesis creation. This approach encourages models to challenge established norms and explore unconventional ideas.

2. **Reducing Historical Bias:** Apply techniques like data augmentation and debiasing Bolukbasi et al. [2016], Liang et al. [2020], Dhamala et al. [2021] to broaden the scope of input datasets, thereby reducing the overfitting of models to historical patterns. By diversifying inputs, models are more likely to generate hypotheses that go beyond traditional paradigms.

3. **Enhancing Human-AI Collaboration:** Develop hybrid systems where LLMs propose unconventional hypotheses, which are refined and validated through iterative feedback from domain experts. This collaborative approach combines AI's ability to explore vast hypothesis spaces with human expertise in assessing practicality and relevance.

4. **Promoting Cross-Domain Knowledge Transfer:** Integrate insights from unrelated fields, such as applying AI methodologies to materials science or borrowing biological frameworks for AI development. The convergence of knowledge from multiple disciplines often results in unexpected connections and innovative hypotheses.

5. **Establishing Novelty Metrics:** Introduce specific metrics to evaluate the novelty of generated hypotheses, such as divergence from existing literature or exploration of uncharted parameter spaces. Explicit metrics ensure that models prioritize innovation and generate hypotheses that challenge existing knowledge.

## 6.2 Improving Feasibility Assessments

Hypothesis validation methods often rely heavily on computational proxies or simulated environments to evaluate feasibility. While these approaches provide scalability and cost efficiency, they frequently fall short of aligning with real-world experimental outcomes. This disconnect is particularly problematic in domains like biomedicine and materials science, where empirical validation is not only preferred but often essential for establishing the reliability of generated hypotheses. The absence of robust empirical validation mechanisms undermines the credibility of computationally derived hypotheses, creating a gap between theoretical predictions and practical applicability.

These challenges carry significant implications. Discrepancies between computational predictions and experimental results erode trust in the reliability of hypotheses, limiting their adoption in critical fields. Domains such as drug development and structural engineering, which demand rigorous empirical evidence, face particular constraints in leveraging computational validation approaches. Furthermore, the lack of alignment with experimental outcomes slows progress in high-stakes applications where feasibility is a prerequisite for real-world deployment. Addressing these limitations is essential for ensuring that hypothesis validation systems fulfill their potential to drive impactful and trustworthy scientific advancements.

**Strategies for Overcoming Limitations:**

1. **Integrating Experimental Validation Pipelines:** Develop frameworks that link computational predictions with empirical validation methods, such as automated laboratory systems or experimental datasets. This integration enhances reliability by directly connecting computational outputs to real-world results.

2. **Establishing Iterative Feedback Mechanisms:** Create feedback loops where experimental outcomes are used to iteratively refine computational models. These mechanisms ensure models continuously adapt to real-world constraints, improving their practical relevance.

3. **Adopting Sim-to-Real Transfer Techniques:** Employ methodologies like domain randomization to bridge the gap between simulated environments and experimental realities. These techniques enhance the robustness and applicability of computational predictions in real-world settings.

4. **Developing Feasibility-Centric Metrics:** Introduce metrics specifically designed to quantify feasibility, such as experimental reproducibility scores or validation success rates. Such metrics provide a systematic way to assess the practical applicability of hypotheses.

5. **Fostering Interdisciplinary Collaborations:** Promote partnerships between computational and experimental researchers to ground hypotheses in practical constraints. These collaborations leverage diverse expertise to enhance the feasibility and applicability of generated hypotheses.

## 6.3 Addressing Domain-Specific and Interdisciplinary Challenges

Hypothesis creation and validation approaches frequently encounter significant challenges in addressing the specialized needs of distinct scientific domains. Each domain often presents unique data structures, terminologies, and methodologies that must be accounted for to ensure accurate and meaningful outcomes. These challenges are further compounded in interdisciplinary research, where integrating diverse knowledge bases and reconciling inconsistent standards across fields is critical. For example, the synthesis of data from biology and artificial intelligence demands a balance between biological intricacies and computational frameworks.

The implications of these challenges are profound. In highly specialized fields like synthetic biology or astrophysics, hypothesis generation approaches may struggle to accommodate domain-specific complexities, thereby reducing their effectiveness. Similarly, interdisciplinary collaboration is often hindered by mismatched methodologies and datasets, slowing progress in areas that rely on cross-domain insights. Emerging fields like bioinformatics or AI-driven materials science exemplify the importance of seamless integration, as the transferability of insights across disciplines is crucial for innovation. Addressing these limitations is essential for advancing hypothesis creation and validation systems that are both domain-sensitive and interdisciplinary in scope.

**Strategies for Overcoming Limitations:**

1. **Domain-Specific Fine-Tuning:** Tailor hypothesis generation models to specific domains using curated datasets and terminologies. This approach enhances model relevance and accuracy by aligning with the standards of the targeted field.

2. **Cross-Disciplinary Ontologies:** Create standardized frameworks and shared ontologies to integrate knowledge across fields. This solution bridges terminological and methodological gaps, enabling seamless interdisciplinary hypothesis generation.

3. **Knowledge Graphs for Interdisciplinary Research:** Leverage knowledge graphs to connect domain-specific concepts and uncover potential synergies between fields. These graphs provide a structured framework for exploring cross-disciplinary relationships.

4. **Collaborative Platforms:** Develop platforms that enable researchers from diverse domains to co-create, refine, and validate hypotheses collaboratively. These platforms facilitate communication and integration of expertise, reducing barriers to interdisciplinary innovation.

5. **Hybrid Models Combining Domain Expertise and General Knowledge:** Integrate domain-specific models with general-purpose LLMs to capitalize on both specialized and broad knowledge. This hybrid approach ensures hypotheses are rigorous and grounded while benefiting from interdisciplinary insights.

## 6.4 Mitigating Data Limitations and Biases

Many hypothesis creation and validation approaches depend heavily on datasets that are often biased, incomplete, or fail to represent the complexities of emerging scientific fields. These shortcomings significantly impact the diversity and generalizability of the generated hypotheses, limiting their applicability across underrepresented areas or demographics. Additionally, biases inherent in training data risk perpetuating systemic inequities, which can undermine trust in AI-driven discoveries and hinder the adoption of these technologies in critical applications.

The implications of these limitations are far-reaching. Hypotheses that are not inclusive of diverse datasets may fail to address the unique challenges of specific populations or regions, limiting their real-world utility. Furthermore, systemic biases introduced during hypothesis creation can erode confidence in the fairness and reliability of scientific discoveries. In emerging fields, where high-quality data is often fragmented or scarce, these issues constrain innovation and slow the progress of groundbreaking research. Addressing these challenges is essential to ensure that hypothesis creation and validation systems promote equity, inclusivity, and the advancement of science in all domains.

**Strategies for Overcoming Limitations:**

1. **Data Augmentation Techniques:** Apply methods such as synthetic data generation, cross-domain transfer, or oversampling to enrich datasets in underrepresented areas. Augmented datasets improve diversity and enable comprehensive hypothesis exploration.

2. **Bias Detection and Mitigation Frameworks:** Implement tools that identify and correct biases, such as fairness metrics or adversarial debiasing. These frameworks ensure equitable and unbiased hypothesis generation.

3. **Collaborative Dataset Curation:** Promote interdisciplinary and community-driven efforts to curate diverse, high-quality datasets. Collaborative curation reduces blind spots and enhances dataset coverage.

4. **Active Learning Approaches:** Use active learning to iteratively refine datasets by prioritizing samples that maximize model performance. This strategy optimizes data collection and improves dataset representativeness.

5. **Openness and Transparency in Dataset Usage:** Advocate for open-access repositories and standardized metadata to detail dataset limitations and biases. Transparency fosters accountability and helps users interpret results in context.

## 6.5 Enhancing Interpretability and Transparency

AI-based hypothesis creation and validation models frequently lack interpretability and transparency, posing significant challenges to their adoption and utility. The intricate nature of these models often obscures the logic and reasoning behind the generated hypotheses, making it difficult for researchers to assess their validity and trustworthiness. This opacity not only limits the understanding of the underlying processes but also creates barriers to effective validation and refinement by domain experts.

The implications of these challenges are profound. A lack of transparency reduces trust in AI-driven hypotheses, particularly in high-stakes areas like healthcare and public policy, where decisions can have far-reaching consequences. Furthermore, domain experts may struggle to engage meaningfully with these systems, as opaque methodologies limit their ability to critically evaluate and refine

hypotheses. This lack of clarity also hinders interdisciplinary collaboration, as researchers from diverse fields may find it difficult to reconcile opaque AI-generated results with established domain-specific practices. Enhancing interpretability and transparency in AI models is therefore essential for fostering trust, enabling effective validation, and promoting collaborative advancements in hypothesis-driven research.

**Strategies for Overcoming Limitations:**

1. **Adoption of Explainable AI (XAI) Techniques:** Incorporate tools like saliency maps, SHAP (SHapley Additive exPlanations), and LIME (Local Interpretable Model-Agnostic Explanations) to clarify AI-generated hypotheses. These techniques improve user understanding and trust.

2. **Human-Readable Outputs:** Design models to output hypotheses in interpretable formats, such as graphs, natural language summaries, or decision trees. Clear outputs facilitate hypothesis validation and enhance communication among researchers.

3. **Model Auditing and Validation:** Establish auditing frameworks to evaluate the transparency of model predictions and hypothesis generation processes. Auditing promotes accountability and identifies inaccuracies in reasoning.

4. **Interactive Visualization Tools:** Develop platforms that allow users to explore model-generated hypotheses and underlying data visually. Visualization tools help identify patterns, anomalies, and areas needing further investigation.

5. **Open-Sourcing Models and Workflows:** Support open-source development of AI models and validation pipelines to enable independent scrutiny and reproducibility. Open access increases transparency and fosters interdisciplinary collaboration.

## 6.6 Encouraging High-Risk, High-Reward Hypotheses

Current frameworks for hypothesis creation and validation frequently emphasize safety and reliability, which, while valuable, often discourage the pursuit of high-risk, high-reward hypotheses. By prioritizing incremental progress and aligning closely with established paradigms, these frameworks inadvertently constrain the exploration of transformative ideas that challenge conventional norms. This cautious approach limits the potential to venture into uncharted territories where groundbreaking discoveries might emerge.

The implications of this risk-averse mindset are significant. It restricts the scientific community's ability to achieve paradigm-shifting advancements, particularly in fields where unconventional thinking is essential. Researchers may be deterred from exploring novel or controversial ideas due to the perceived risks and lack of institutional support, further reinforcing existing knowledge silos. This focus on low-risk, incremental progress slows innovation, particularly in emerging fields where disruptive breakthroughs are most needed. Developing frameworks that balance reliability with the freedom to explore bold hypotheses is critical for driving transformative scientific and technological advancements.

**Strategies for Overcoming Limitations:**

1. **Incorporating Risk-Tolerant Metrics:** Introduce evaluation metrics that reward hypotheses with high potential impact despite their inherent uncertainty. This approach fosters exploration of transformative ideas while maintaining a balance with feasibility.

2. **Exploratory Funding Models:** Advocate for funding mechanisms, such as "moonshot" grants, to support high-risk, high-reward research. Dedicated financial support reduces barriers and incentivizes bold innovations.

3. **Scenario-Based Validation Frameworks:** Create frameworks that assess hypotheses across diverse scenarios to evaluate long-term risks and rewards. Scenario-based evaluations provide a nuanced understanding of high-risk ideas.

4. **Cross-Disciplinary Teams for Validation:** Involve interdisciplinary teams to evaluate unconventional hypotheses from multiple perspectives. Diverse expertise ensures a balanced assessment of high-risk proposals.

5. **Recognition and Incentive Structures:** Establish systems that reward researchers for pursuing innovative ideas, regardless of immediate success. Positive incentives promote a culture of creativity and exploration.

## 6.7 Scalability in Data Integration

Hypothesis generation and validation approaches often encounter significant scalability challenges when tasked with integrating large-scale, heterogeneous datasets. As the volume and diversity of available data continue to grow, computational demands increase exponentially, complicating the effective processing and analysis of these datasets. The complexity of managing such multidimensional, real-world data further exacerbates the difficulty of deriving meaningful insights.

These scalability limitations have notable implications for scientific discovery. They hinder the ability to generate insights that leverage information from multiple domains or large-scale datasets, reducing the potential for cross-disciplinary innovation. Additionally, the increased computational costs and decreased efficiency associated with handling vast datasets strain hypothesis workflows, making them less practical for widespread application. Without addressing these challenges, current approaches risk becoming obsolete in the face of the ever-expanding data landscape, underscoring the need for more scalable and efficient solutions.

**Strategies for Overcoming Limitations:**

1. **Scalable Data Processing Architectures:** Utilize distributed computing frameworks, such as Apache Spark or TensorFlow Extended, to handle large-scale datasets. Parallel processing reduces bottlenecks in data-intensive tasks.

2. **Federated Learning for Data Integration:** Implement federated learning to integrate distributed data sources without centralizing them. This preserves privacy while enabling scalable analysis across systems.

3. **Advanced Data Preprocessing Pipelines:** Develop automated workflows for normalizing, cleaning, and aligning heterogeneous datasets. Preprocessing ensures data quality and consistency for improved hypothesis generation.

4. **Knowledge Graphs for Data Fusion:** Use knowledge graphs to unify datasets by representing relationships across sources. Semantic linking facilitates the discovery of hidden insights and relationships.

5. **Optimized Query and Storage Systems:** Deploy systems like graph databases or NoSQL solutions to manage multidimensional data efficiently. These tools enhance storage and retrieval capabilities, improving real-time integration.

## 6.8 Feedback-Driven Refinement

Hypothesis creation and validation systems frequently lack robust feedback mechanisms that enable the integration of iterative expert input. This absence of interaction diminishes the domain-specific relevance of generated hypotheses, as these systems are unable to benefit from nuanced corrections or refinements provided by subject matter experts. Furthermore, without continuous feedback loops, these models struggle to adapt effectively to evolving research priorities or new datasets, limiting their overall utility.

The implications of this limitation are significant. The exclusion of expert insights reduces the quality and contextual accuracy of hypotheses, often leading to outputs that are less applicable or meaningful in real-world scenarios. Additionally, the inability to identify and correct errors in AI-generated hypotheses undermines trust in these systems and restricts their adoption in critical fields. To remain relevant and impactful, hypothesis creation and validation frameworks must prioritize the development of adaptable, expert-informed feedback mechanisms that align with dynamic research needs.

**Strategies for Overcoming Limitations:**

1. **Human-in-the-Loop Systems:** Design systems that allow experts to refine hypotheses iteratively through feedback loops. This ensures hypotheses align with domain-specific standards and relevance.

2. **Interactive Model Interfaces:** Create user-friendly interfaces for visualizing, modifying, and providing feedback on hypotheses. Interactive tools enhance transparency and facilitate dynamic refinement.

3. **Reinforcement Learning with Human Feedback:** Use reinforcement learning algorithms to optimize models based on expert feedback. Continuous improvement adapts models to evolving research needs.

4. **Iterative Validation Pipelines:** Develop pipelines where hypotheses undergo cycles of AI validation and expert evaluation. Iterative approaches improve robustness and minimize errors.

5. **Crowdsourced Feedback Mechanisms:** Incorporate diverse insights through crowdsourcing platforms for collaborative hypothesis refinement. Crowdsourcing broadens perspectives and enhances inclusivity in research processes.

Addressing limitations in hypothesis creation and validation is crucial for realizing the full potential of LLMs in scientific discovery. By fostering novelty, improving feasibility, and enhancing interpretability, these tools can become robust enablers of transformative research. Tackling challenges like data biases, scalability, and interdisciplinary integration will empower researchers to drive breakthroughs across diverse fields. Strategies such as interdisciplinary frameworks, human-in-the-loop systems, and novel evaluation metrics are essential for encouraging bold, high-impact hypotheses. Implementing these advancements will bridge existing gaps and position LLMs as indispensable allies in accelerating innovation and fostering a new era of scientific discovery.

# 7 Conclusion

This survey has explored the transformative role of Large Language Models (LLMs) in scientific hypothesis creation and validation, examining their methodologies, datasets, limitations, and future directions. LLMs exhibit exceptional capabilities in processing vast datasets, identifying patterns, and generating insights, leveraging approaches such as knowledge graphs, retrieval-augmented generation, predictive modeling, and simulation-based validation. These methods offer diverse strengths, yet they remain constrained by challenges such as limited novelty, data biases, domain-specific constraints, and the reliance on computational proxies rather than empirical validation. To fully harness the potential of LLMs in scientific discovery, future research must focus on advancing novelty, feasibility, and interdisciplinary integration. Enhancing novelty through advanced generative exploration frameworks, such as generative adversarial models and contrastive learning, can help overcome the tendency of LLMs to rely on existing knowledge. Strengthening feasibility via hybrid computational-empirical pipelines will ensure that hypotheses generated by LLMs align with experimental validation. Additionally, addressing interdisciplinary challenges through cross-domain ontologies, federated learning, and collaborative AI-human interfaces will enhance the adaptability of LLMs in complex research landscapes.

Ensuring data integrity and mitigating biases is another critical step. Techniques including strategic knowledge graph augmentation, adversarial debiasing methods, and fairness-aware machine learning models can help mitigate biases and enhance the scalability, inclusivity, and trustworthiness of LLM-driven hypothesis generation. Moreover, interpretability-focused systems, incorporating explainable AI techniques and uncertainty quantification, will be essential for bridging the gap between computational insights and human intuition. The future of LLM-driven scientific discovery lies in iterative human-AI collaboration, where researchers and AI systems co-develop hypotheses in a risk-tolerant and adaptive framework. By integrating explainability, risk assessment, and interdisciplinary synthesis, LLMs can serve as powerful catalysts for high-risk, high-reward research, accelerating breakthroughs in biomedicine, materials science, climate modeling, and beyond. In conclusion, as advancements in LLM technology continue to evolve, their role in hypothesis-driven research will expand, redefining scientific exploration. By addressing existing limitations and fostering a collaborative, interpretability-driven ecosystem, LLMs hold the potential to usher in a new era of transformative knowledge creation, interdisciplinary discovery, and scalable scientific innovation.

# References

Gary F Bradshaw, Patrick W Langley, and Herbert A Simon. Studying scientific discovery by computer simulation. *Science*, 222(4627):971–975, 1983.

Herbert A Simon. Scientific discovery as problem solving: Reply to critics. 1992.

Pat Langley and Randolph Jones. A computational model of scientific insight. *The nature of creativity: Contemporary psychological perspectives*, 177(201):2, 1988.

Pat Langley. The computer-aided discovery of scientific knowledge. In *International Conference on Discovery Science*, pages 25–39. Springer, 1998.

Pat Langley. The computational support of scientific discovery. *International Journal of Human-Computer Studies*, 53(3):393–410, 2000.

Sašo Džeroski, Pat Langley, and Ljupčo Todorovski. Computational discovery of scientific knowledge. In *Computational discovery of scientific knowledge: Introduction, techniques, and applications in environmental and life sciences*, pages 1–14. Springer.

Pat Langley and Herbert A Simon. The central role of learning in cognition. In *Cognitive skills and their acquisition*, pages 361–380. Psychology Press, 2013.

Pat Langley. Integrated systems for computational scientific discovery. In *Proceedings of the Thirty-Eighth AAAI Conference on Artificial Intelligence and Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence and Fourteenth Symposium on Educational Advances in Artificial Intelligence*, pages 22598–22606, 2024.

Yoshua Bengio and Yann LeCun. Scaling learning algorithms towards AI. In *Large Scale Kernel Machines*. MIT Press, 2007.

Geoffrey E. Hinton, Simon Osindero, and Yee Whye Teh. A fast learning algorithm for deep belief nets. *Neural Computation*, 18:1527–1554, 2006.

Margaret A. Boden. Creativity and artificial intelligence. *Artificial Intelligence*, 103(1):347–356, 1998. ISSN 0004-3702. doi: https://doi.org/10.1016/S0004-3702(98)00055-1. URL `https://www.sciencedirect.com/science/article/pii/S0004370298000551`. Artificial Intelligence 40 years later.

AI@Meta. Llama 3 model card. 2024. URL `https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md`.

Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024.

Iz Beltagy, Kyle Lo, and Arman Cohan. Scibert: A pretrained language model for scientific text. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3615–3620, 2019.

Biqing Qi, Kaiyan Zhang, Kai Tian, Haoxiang Li, Zhang-Ren Chen, Sihang Zeng, Ermo Hua, Hu Jinfang, and Bowen Zhou. Large language models as biomedical hypothesis generators: A comprehensive evaluation. *arXiv preprint arXiv:2407.08940*, 2024.

Qingyun Wang, Doug Downey, Heng Ji, and Tom Hope. Scimon: Scientific inspiration machines optimized for novelty. *arXiv preprint arXiv:2305.14259*, 2023.

Guanjie Wang, Jingjing Hu, Jian Zhou, Sen Liu, Qingjiang Li, and Zhimei Sun. Knowledge-guided large language model for material science. *Review of Materials Research*, page 100007, 2025. ISSN 3050-9130. doi: https://doi.org/10.1016/j.revmat.2025.100007. URL `https://www.sciencedirect.com/science/article/pii/S3050913025000075`.

Rui Ding, Jianguo Liu, Kang Hua, Xuebin Wang, Xiaoben Zhang, Minhua Shao, Yuxin Chen, and Junhong Chen. Leveraging data mining, active learning, and domain adaptation for efficient discovery of advanced oxygen evolution electrocatalysts. *Science Advances*, 11(14):eadr9038, 2025. doi: 10.1126/sciadv.adr9038. URL https://www.science.org/doi/abs/10.1126/sciadv.adr9038.

Rong Ji, Kai Gong, Lihong Huang, Wenxian Yang, and Rongshan Yu. Leveraging llms for automated analysis of biomedical data. In *2024 9th International Conference on Communication, Image and Signal Processing (CCISP)*, pages 67–71, 2024. doi: 10.1109/CCISP63826.2024.10765518.

Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1. MIT Press, 2016.

Justin Sybrandt, Micheal Shtutman, and Ilya Safro. Large-scale validation of hypothesis generation systems via candidate ranking. In *2018 IEEE International Conference on Big Data (Big Data)*, pages 1494–1503. IEEE, 2018.

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

Dian-Zhao Lin, Kai-Jui Pan, Yuyin Li, Charles B. null, Lingyu Zhang, Krish N. Jayarapu, Tianchen Li, Jasmine Vy Tran, William A. Goddard, Zhengtang Luo, and Yayuan Liu. A high-throughput experimentation platform for data-driven discovery in electrochemistry. *Science Advances*, 11(14): eadu4391, 2025. doi: 10.1126/sciadv.adu4391. URL https://www.science.org/doi/abs/10.1126/sciadv.adu4391.

Andres M Bran, Sam Cox, Oliver Schilter, Carlo Baldassari, Andrew D White, and Philippe Schwaller. Chemcrow: Augmenting large-language models with chemistry tools. *arXiv preprint arXiv:2304.05376*, 2023.

Muamer Kadic, Graeme W Milton, Martin van Hecke, and Martin Wegener. 3d metamaterials. *Nature Reviews Physics*, 1(3):198–210, 2019.

Katia Bertoldi, Vincenzo Vitelli, Johan Christensen, and Martin Van Hecke. Flexible mechanical metamaterials. *Nature Reviews Materials*, 2(11):1–11, 2017.

Zian Jia, Fan Liu, Xihang Jiang, and Lifeng Wang. Engineering lattice metamaterials for extreme property, programmability, and multifunctionality. *Journal of Applied Physics*, 127(15), 2020.

Pengcheng Jiao and Amir H Alavi. Artificial intelligence-enabled smart mechanical metamaterials: advent and future trends. *International Materials Reviews*, 66(6):365–393, 2021.

Jens Bauer, Lucas R Meza, Tobias A Schaedler, Ruth Schwaiger, Xiaoyu Zheng, and Lorenzo Valdevit. Nanolattices: an emerging class of mechanical metamaterials. *Advanced Materials*, 29 (40):1701850, 2017.

Ioannis Papadimitriou, Ilias Gialampoukidis, Stefanos Vrochidis, and Ioannis Kompatsiaris. Ai methods in materials design, discovery and manufacturing: A review. *Computational Materials Science*, 235:112793, 2024.

Zonglin Yang, Xinya Du, Junxian Li, Jie Zheng, Soujanya Poria, and Erik Cambria. Large language models for automated open-domain scientific hypotheses discovery. *arXiv preprint arXiv:2309.02726*, 2023.

Parshin Shojaee, Kazem Meidani, Shashank Gupta, Amir Barati Farimani, and Chandan K Reddy. Llm-sr: Scientific equation discovery via programming with large language models. *arXiv preprint arXiv:2404.18400*, 2024.

Bernardino Romera-Paredes, Mohammadamin Barekatain, Alexander Novikov, Matej Balog, M Pawan Kumar, Emilien Dupont, Francisco JR Ruiz, Jordan S Ellenberg, Pengming Wang, Omar Fawzi, et al. Mathematical discoveries from program search with large language models. *Nature*, 625(7995):468–475, 2024.

Trieu H Trinh, Yuhuai Wu, Quoc V Le, He He, and Thang Luong. Solving olympiad geometry without human demonstrations. *Nature*, 625(7995):476–482, 2024.

Nicholas A Lesica, Nishchay Mehta, Joseph G Manjaly, Li Deng, Blake S Wilson, and Fan-Gang Zeng. Harnessing the power of artificial intelligence to transform hearing healthcare and research. *Nature Machine Intelligence*, 3(10):840–849, 2021.

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

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

K Chen et al. Chemist-x: Large language model-empowered agent for reaction condition recommendation in chemical synthesis. *Preprint at*, 2024.

Yangqiaoyu Zhou, Haokun Liu, Tejes Srivastava, Hongyuan Mei, and Chenhao Tan. Hypothesis generation with large language models, 2024. URL https://arxiv.org/abs/2404.04326.

Justin Sybrandt, Michael Shtutman, and Ilya Safro. Moliere: Automatic biomedical hypothesis generation system. In *Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1633–1642, 2017.

Kishlay Jha, Guangxu Xun, Yaqing Wang, and Aidong Zhang. Hypothesis generation from text based on co-evolution of biomedical concepts. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 843–851, 2019.

Alhussein Fawzi, Matej Balog, Aja Huang, Thomas Hubert, Bernardino Romera-Paredes, Mohammadamin Barekatain, Alexander Novikov, Francisco J R Ruiz, Julian Schrittwieser, Grzegorz Swirszcz, et al. Discovering faster matrix multiplication algorithms with reinforcement learning. *Nature*, 610(7930):47–53, 2022.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

Jingyuan Qi, Zhiyang Xu, Ying Shen, Minqian Liu, Di Jin, Qifan Wang, and Lifu Huang. The art of SOCRATIC QUESTIONING: Recursive thinking with large language models. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 4177–4199, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.255. URL https://aclanthology.org/2023.emnlp-main.255.

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

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

Anubhav Jain, Joseph Montoya, Shyam Dwaraknath, Nils ER Zimmermann, John Dagdelen, Matthew Horton, Patrick Huck, Donny Winston, Shreyas Cholia, Shyue Ping Ong, et al. The materials project: Accelerating materials design through theory-driven data and tools. *Handbook of Materials Modeling: Methods: Theory and Modeling*, pages 1751–1784, 2020.

Luis Tari, Saadat Anwar, Shanshan Liang, James Cai, and Chitta Baral. Discovering drug–drug interactions: a text-mining and reasoning approach based on properties of drug metabolism. *Bioinformatics*, 26(18):i547–i553, 2010.

Raymond Fok and Daniel S Weld. In search of verifiability: Explanations rarely enable complementary performance in ai-advised decision making. *AI Magazine*, 45(3):317–332, 2024.

Aditya Paul, Chi Lok Yu, Eva Adelina Susanto, Nicholas Wai Long Lau, and Gwenyth Isobel Meadows. Agentpeertalk: Empowering students through agentic-ai-driven discernment of bullying and joking in peer interactions in schools, 2024. URL `https://arxiv.org/abs/2408.01459`.

Jules White. Building living software systems with generative & agentic ai, 2024. URL `https://arxiv.org/abs/2408.01768`.

Alan Chan, Rebecca Salganik, Alva Markelius, Chris Pang, Nitarshan Rajkumar, Dmitrii Krasheninnikov, Lauro Langosco, Zhonghao He, Yawen Duan, Micah Carroll, Michelle Lin, Alex Mayhew, Katherine Collins, Maryam Molamohammadi, John Burden, Wanru Zhao, Shalaleh Rismani, Konstantinos Voudouris, Umang Bhatt, Adrian Weller, David Krueger, and Tegan Maharaj. Harms from increasingly agentic algorithmic systems. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '23, page 651–666, New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9798400701924. doi: 10.1145/3593013.3594033. URL `https://doi.org/10.1145/3593013.3594033`.

Antonin Sulc, Thorsten Hellert, Raimund Kammering, Hayden Houscher, and Jason St. John. Towards agentic ai on particle accelerators, 2024. URL `https://arxiv.org/abs/2409.06336`.

Huachuan Qiu and Zhenzhong Lan. Interactive agents: Simulating counselor-client psychological counseling via role-playing llm-to-llm interactions. *arXiv preprint arXiv:2408.15787*, 2024.

Mourad Gridach, Jay Nanavati, Khaldoun Zine El Abidine, Lenon Mendes, and Christina Mack. Agentic ai for scientific discovery: A survey of progress, challenges, and future directions. *arXiv preprint arXiv:2503.08979*, 2025.

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

Nate Gruver, Anuroop Sriram, Andrea Madotto, Andrew Gordon Wilson, C Lawrence Zitnick, and Zachary Ulissi. Fine-tuned language models generate stable inorganic materials as text. *arXiv preprint arXiv:2402.04379*, 2024.

Xiangru Tang, Qiao Jin, Kunlun Zhu, Tongxin Yuan, Yichi Zhang, Wangchunshu Zhou, Meng Qu, Yilun Zhao, Jian Tang, Zhuosheng Zhang, et al. Prioritizing safeguarding over autonomy: Risks of llm agents for science. *arXiv preprint arXiv:2402.04247*, 2024.

Yonadav Shavit, Sandhini Agarwal, Miles Brundage, Steven Adler, Cullen O'Keefe, Rosie Campbell, Teddy Lee, Pamela Mishkin, Tyna Eloundou, Alan Hickey, et al. Practices for governing agentic ai systems. *Research Paper, OpenAI, December*, 2023.

Yu Zhang, Xiusi Chen, Bowen Jin, Sheng Wang, Shuiwang Ji, Wei Wang, and Jiawei Han. A comprehensive survey of scientific large language models and their applications in scientific discovery. *arXiv preprint arXiv:2406.10833*, 2024a.

Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson, Shyamal Buch, Dallas Card, Rodrigo Castellon, Niladri Chatterji, Annie Chen, Kathleen Creel, Jared Quincy Davis, Dora Demszky, Chris Donahue, Moussa Doumbouya, Esin Durmus, Stefano Ermon, John Etchemendy, Kawin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale, Lauren Gillespie, Karan Goel, Noah Goodman, Shelby Grossman, Neel Guha, Tatsunori Hashimoto, Peter Henderson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang, Thomas Icard, Saahil Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keeling, Fereshte Khani, Omar Khattab, Pang Wei Koh, Mark Krass, Ranjay Krishna, Rohith Kuditipudi, Ananya Kumar, Faisal Ladhak, Mina Lee, Tony Lee, Jure Leskovec, Isabelle Levent, Xiang Lisa Li, Xuechen Li, Tengyu Ma, Ali Malik, Christopher D. Manning, Suvir Mirchandani, Eric Mitchell, Zanele Munyikwa, Suraj Nair, Avanika Narayan, Deepak Narayanan, Ben Newman, Allen Nie, Juan Carlos Niebles, Hamed Nilforoshan, Julian Nyarko, Giray Ogut, Laurel Orr, Isabel Papadimitriou, Joon Sung Park, Chris Piech, Eva Portelance, Christopher Potts, Aditi Raghunathan, Rob Reich, Hongyu Ren, Frieda Rong, Yusuf Roohani, Camilo Ruiz, Jack Ryan, Christopher Ré, Dorsa Sadigh, Shiori Sagawa, Keshav Santhanam, Andy Shih, Krishnan Srinivasan, Alex Tamkin, Rohan Taori,

Armin W. Thomas, Florian Tramèr, Rose E. Wang, William Wang, Bohan Wu, Jiajun Wu, Yuhuai Wu, Sang Michael Xie, Michihiro Yasunaga, Jiaxuan You, Matei Zaharia, Michael Zhang, Tianyi Zhang, Xikun Zhang, Yuhui Zhang, Lucia Zheng, Kaitlyn Zhou, and Percy Liang. On the opportunities and risks of foundation models, 2021.

Raphael Schumann, Wanrong Zhu, Weixi Feng, Tsu-Jui Fu, Stefan Riezler, and William Yang Wang. Velma: Verbalization embodiment of llm agents for vision and language navigation in street view. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 18924–18933, 2024.

Manpreet S Katari, Steve D Nowicki, Felipe F Aceituno, Damion Nero, Jonathan Kelfer, Lee Parnell Thompson, Juan M Cabello, Rebecca S Davidson, Arthur P Goldberg, Dennis E Shasha, et al. Virtualplant: a software platform to support systems biology research. *Plant physiology*, 152(2): 500–515, 2010.

Oskar Wysocki, Magdalena Wysocka, Danilo Carvalho, Alex Teodor Bogatu, Danilo Miranda Gusicuma, Maxime Delmas, Harriet Unsworth, and Andre Freitas. An llm-based knowledge synthesis and scientific reasoning framework for biomedical discovery. *arXiv preprint arXiv:2406.18626*, 2024.

Alexandre Blanco-Gonzalez, Alfonso Cabezon, Alejandro Seco-Gonzalez, Daniel Conde-Torres, Paula Antelo-Riveiro, Angel Pineiro, and Rebeca Garcia-Fandino. The role of ai in drug discovery: challenges, opportunities, and strategies. *Pharmaceuticals*, 16(6):891, 2023.

Ilya Tyagin and Ilya Safro. Dyport: dynamic importance-based biomedical hypothesis generation benchmarking technique. *BMC bioinformatics*, 25, 2024.

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

Samuel Schmidgall, Rojin Ziaei, Carl Harris, Eduardo Reis, Jeffrey Jopling, and Michael Moor. Agentclinic: a multimodal agent benchmark to evaluate ai in simulated clinical environments. *arXiv preprint arXiv:2405.07960*, 2024.

Jung-Hun Kim and Aviv Segev. Research hypothesis generation using link prediction in a bipartite graph. In *2018 IEEE International Conference on Big Data (Big Data)*, pages 2863–2867. IEEE, 2018.

Yubo Ma, Zhibin Gou, Junheng Hao, Ruochen Xu, Shuohang Wang, Liangming Pan, Yujiu Yang, Yixin Cao, Aixin Sun, Hany Awadalla, et al. Sciagent: Tool-augmented language models for scientific reasoning. *arXiv preprint arXiv:2402.11451*, 2024.

Marianne Aubin Le Quéré, Hope Schroeder, Casey Randazzo, Jie Gao, Ziv Epstein, Simon Tangi Perrault, David Mimno, Louise Barkhuus, and Hanlin Li. Llms as research tools: Applications and evaluations in hci data work. In *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems*, pages 1–7, 2024.

National Center for Biotechnology Information. PubMed: A resource for biomedical literature, 2025. URL https://pubmed.ncbi.nlm.nih.gov/.

Michael Ashburner, Catherine A Ball, Judith A Blake, David Botstein, Heather Butler, J Michael Cherry, Allan P Davis, Kara Dolinski, Selina S Dwight, Janan T Eppig, et al. Gene ontology: tool for the unification of biology. *Nature genetics*, 25(1):25–29, 2000.

Megan C Conroy, Ben Lacey, Jelena Bešević, Wemimo Omiyale, Qi Feng, Mark Effingham, Jonathan Sellers, Simon Sheard, Mahesh Pancholi, Gareth Gregory, et al. Uk biobank: a globally important resource for cancer research. *British Journal of Cancer*, 128(4):519–527, 2023.

Theo Walker, Christopher M Grulke, Diane Pozefsky, and Alexander Tropsha. Chembench: a cheminformatics workbench. *Bioinformatics*, 26(23):3000–3001, 2010.

Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, Zac Kenton, Sasha Brown, Will Hawkins, Tom Stepleton, Courtney Biles, Abeba Birhane, Julia Haas, Laura Rimell, Lisa Anne Hendricks, William Isaac, Sean Legassick, Geoffrey Irving, and Iason Gabriel. Ethical and social risks of harm from language models, 2021.

Haigen Hu, Xiaoyuan Wang, Yan Zhang, Qi Chen, and Qiu Guan. A comprehensive survey on contrastive learning. *Neurocomputing*, page 128645, 2024.

Jiawei Han and Yongjian Fu. Dynamic generation and refinement of concept hierarchies for knowledge discovery in databases. In *KDD workshop*, pages 157–168, 1994.

Marianne Aubin Le Quéré, Hope Schroeder, Casey Randazzo, Jie Gao, Ziv Epstein, Simon Tangi Perrault, David Mimno, Louise Barkhuus, and Hanlin Li. Llms as research tools: Applications and evaluations in hci data work. In *Extended Abstracts of the 2024 CHI Conference on Human Factors in Computing Systems*, CHI EA '24, New York, NY, USA, 2024. Association for Computing Machinery. ISBN 9798400703317. doi: 10.1145/3613905.3636301. URL https://doi.org/10.1145/3613905.3636301.

Ross D King. Rise of the robo scientists. *Scientific American*, 304(1):72–77, 2011.

Hatem Fakhruldeen, Gabriella Pizzuto, Jakub Glowacki, and Andrew Ian Cooper. Archemist: Autonomous robotic chemistry system architecture. In *2022 International Conference on Robotics and Automation (ICRA)*, pages 6013–6019. IEEE, 2022.

Kevin Williams, Elizabeth Bilsland, Andrew Sparkes, Wayne Aubrey, Michael Young, Larisa N Soldatova, Kurt De Grave, Jan Ramon, Michaela De Clare, Worachart Sirawaraporn, et al. Eve: Integration of machine learning with compound testing in a robot scientist. In *Antenna Live: Robot Scientist, Location: London*. Science Museum, 2015.

Pranav Narayanan Venkit, Tatiana Chakravorti, Vipul Gupta, Heidi Biggs, Mukund Srinath, Koustava Goswami, Sarah Rajtmajer, and Shomir Wilson. An audit on the perspectives and challenges of hallucinations in nlp. *arXiv preprint arXiv:2404.07461*, 2024.

Mohamad Yaser Jaradeh, Allard Oelen, Kheir Eddine Farfar, Manuel Prinz, Jennifer D'Souza, Gábor Kismihók, Markus Stocker, and Sören Auer. Open research knowledge graph: next generation infrastructure for semantic scholarly knowledge. In *Proceedings of the 10th international conference on knowledge capture*, pages 243–246, 2019.

DeBrae Kennedy-Mayo and Jake Gord. "model cards for model reporting" in 2024: Reclassifying category of ethical considerations in terms of trustworthiness and risk management. In *Future of Information and Communication Conference*, pages 179–196. Springer, 2025.

José Luiz Nunes, Gabriel DJ Barbosa, Clarisse Sieckenius de Souza, and Simone DJ Barbosa. Using model cards for ethical reflection on machine learning models: an interview-based study. *Journal on Interactive Systems*, 15(1):1–19, 2024.

Ulrich Witt. Propositions about novelty. *Journal of Economic Behavior & Organization*, 70(1-2):311–320, 2009.

John E Hallsworth, Zulema Udaondo, Carlos Pedrós-Alió, Juan Höfer, Kathleen C Benison, Karen G Lloyd, Radamés JB Cordero, Claudia BL de Campos, Michail M Yakimov, and Ricardo Amils. Scientific novelty beyond the experiment. *Microbial Biotechnology*, 16(6):1131–1173, 2023.

Sotaro Shibayama, Yutaro Baba, and John P Walsh. Measuring novelty in science with word embedding. *PloS one*, 16(7):e0254034, 2021. doi: 10.1371/journal.pone.0254034.

Jaeill Kim, Duhun Hwang, Eunjung Lee, Jangwon Suh, Jimyeong Kim, and Wonjong Rhee. Enhancing contrastive learning with efficient combinatorial positive pairing, 2024. URL https://arxiv.org/abs/2401.05730.

Michael G Walker. How feasible is automated discovery? *IEEE Intelligent Systems*, 2(01):69–82, 1987.

Junyue Song, Xin Wu, and Yi Cai. Step feasibility-aware and error-correctable entailment tree generation. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 15296–15308, 2024.

Chong Wang, Mengyao Li, Junjun He, Zhongruo Wang, Erfan Darzi, Zan Chen, Jin Ye, Tianbin Li, Yanzhou Su, Jing Ke, et al. A survey for large language models in biomedicine. *arXiv preprint arXiv:2409.00133*, 2024a.

U.S. National Library of Medicine. Medical Subject Headings (MeSH), 2025. URL `https://www.nlm.nih.gov/mesh/meshhome.html`.

Barbara Zdrazil, Eloy Felix, Fiona Hunter, Emma J Manners, James Blackshaw, Sybilla Corbett, Marleen de Veij, Harris Ioannidis, David Mendez Lopez, Juan F Mosquera, Maria Paula Magarinos, Nicolas Bosc, Ricardo Arcila, Tevfik Kizilören, Anna Gaulton, A Patrícia Bento, Melissa F Adasme, Peter Monecke, Gregory A Landrum, and Andrew R Leach. The chembl database in 2023: a drug discovery platform spanning multiple bioactivity data types and time periods. *Nucleic Acids Research*, 52(D1):D1180–D1192, 11 2023. ISSN 0305-1048. doi: 10.1093/nar/gkad1004. URL `https://doi.org/10.1093/nar/gkad1004`.

J-D Kim, Tomoko Ohta, Yuka Tateisi, and Jun'ichi Tsujii. Genia corpus—a semantically annotated corpus for bio-textmining. *Bioinformatics*, 19(suppl_1):i180–i182, 2003.

Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta, and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs. *Advances in neural information processing systems*, 33:22118–22133, 2020.

Alexander Dunn, Qi Wang, Alex Ganose, Daniel Dopp, and Anubhav Jain. Benchmarking materials property prediction methods: the matbench test set and automatminer reference algorithm. *npj Computational Materials*, 6(1):138, 2020.

Prabhat, Karthik Kashinath, Mayur Mudigonda, Sol Kim, Lukas Kapp-Schwoerer, Andre Graubner, Ege Karaismailoglu, Leo von Kleist, Thorsten Kurth, Annette Greiner, et al. Climatenet: An expert-labelled open dataset and deep learning architecture for enabling high-precision analyses of extreme weather. *Geoscientific Model Development Discussions*, 2020:1–28, 2020.

Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer vision–ECCV 2014: 13th European conference, zurich, Switzerland, September 6-12, 2014, proceedings, part v 13*, pages 740–755. Springer, 2014.

Agustín Borrego, Danilo Dessì, Daniel Ayala, Inma Hernández, Francesco Osborne, Diego Reforgiato Recupero, Davide Buscaldi, David Ruiz, and Enrico Motta. Research hypothesis generation over scientific knowledge graphs. *Knowledge-Based Systems*, 315:113280, 2025. ISSN 0950-7051. doi: https://doi.org/10.1016/j.knosys.2025.113280. URL `https://www.sciencedirect.com/science/article/pii/S0950705125003272`.

David S Wishart, Yannick D Feunang, An C Guo, Elvis J Lo, Ana Marcu, Jason R Grant, Tanvir Sajed, Daniel Johnson, Carin Li, Zinat Sayeeda, et al. Drugbank 5.0: a major update to the drugbank database for 2018. *Nucleic acids research*, 46(D1):D1074–D1082, 2018.

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*, 2018.

Anubhav Jain, Shyue Ping Ong, Geoffroy Hautier, Wei Chen, William Davidson Richards, Stephen Dacek, Shreyas Cholia, Dan Gunter, David Skinner, Gerbrand Ceder, et al. Commentary: The materials project: A materials genome approach to accelerating materials innovation. *APL materials*, 1(1), 2013.

Minoru Kanehisa and Susumu Goto. Kegg: kyoto encyclopedia of genes and genomes. *Nucleic acids research*, 28(1):27–30, 2000.

United States Census Bureau. American community survey, 2025. URL `https://www.census.gov/programs-surveys/acs`.

United States Patent and Trademark Office. United states patent and trademark office patent database, 2025. URL `https://www.uspto.gov/patents`.

Shashi Narayan, Shay B Cohen, and Mirella Lapata. Don't give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization. *arXiv preprint arXiv:1808.08745*, 2018.

Sally Bamford, Emily Dawson, Simon Forbes, Jody Clements, Roger Pettett, Ahmet Dogan, A Flanagan, Jon Teague, P Andrew Futreal, Michael R Stratton, et al. The cosmic (catalogue of somatic mutations in cancer) database and website. *British journal of cancer*, 91(2):355–358, 2004.

Rong Wang, Kun Sun, and Jonas Kuhn. Dspy-based neural-symbolic pipeline to enhance spatial reasoning in llms, 2024b. URL `https://arxiv.org/abs/2411.18564`.

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *Advances in neural information processing systems*, 30, 2017.

Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Model-agnostic interpretability of machine learning. *arXiv preprint arXiv:1606.05386*, 2016.

Shuo Ren, Pu Jian, Zhenjiang Ren, Chunlin Leng, Can Xie, and Jiajun Zhang. Towards scientific intelligence: A survey of llm-based scientific agents, 2025. URL `https://arxiv.org/abs/2503.24047`.

TianZhixi Yin, Ruozhu Feng, Jie Bao, Peiyuan Gao, Yangang Liang, Job Heather, Alan Aspuru-Guzik, and Wei Wang. Learning advance: Robotics-llm guided hypotheses generation for the discovery of chemical knowledge. *ChemRxiv*, 2025. doi: 10.26434/chemrxiv-2025-n1b4l.

Mathilde Resell, Elisabeth Pimpisa Graarud, Hanne-Line Rabben, Animesh Sharma, Lars Hagen, Linh Hoang, Nan T. Skogaker, Anne Aarvik, Magnus K. Svensson, Manoj Amrutkar, Caroline S. Verbeke, Surinder K. Batra, Gunnar Qvigstad, Timothy C. Wang, Anil Rustgi, Duan Chen, and Chun-Mei Zhao. Knowledge discovery in datasets of proteomics by systems modeling in translational research on pancreatic cancer. *bioRxiv*, 2025. doi: 10.1101/2025.02.23.639474. URL `https://www.biorxiv.org/content/early/2025/02/27/2025.02.23.639474`.

Zhenyu Yu. Ai for science: A comprehensive review on innovations, challenges, and future directions. *International Journal of Artificial Intelligence for Science (IJAI4S)*, 1(1), 2025.

Thomas Steinecker, Thorsten Luettel, and Mirko Maehlisch. Towards safety aware ai agents. `https://www.researchgate.net/publication/389351017_Towards_Safety_Aware_AI_Agents`, 2025. Preprint on ResearchGate, accessed April 3, 2025.

Amanda Kau, Xuzeng He, Aishwarya Nambissan, Aland Astudillo, Hui Yin, and Amir Aryani. Combining knowledge graphs and large language models. *arXiv preprint arXiv:2407.06564*, 2024.

Jie Bai, Sebastian Mosbach, Christopher J. Taylor, et al. A dynamic knowledge graph approach to distributed self-driving laboratories. *Nature Communications*, 15(1):462, 2024a. doi: 10.1038/s41467-023-44599-9. URL `https://doi.org/10.1038/s41467-023-44599-9`.

Yu Zhang, Zhiyong Cheng, Fan Liu, Xun Yang, and Yuxin Peng. Decoupled domain-specific and domain-conditional representation learning for cross-domain recommendation. *Information Processing & Management*, 61(3):103689, 2024b.

Jiaxin Bai, Yicheng Wang, Tianshi Zheng, Yue Guo, Xin Liu, and Yangqiu Song. Advancing abductive reasoning in knowledge graphs through complex logical hypothesis generation. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1312–1329, 2024b.

C-H Liu, K-L Chang, Jason J-Y Chen, and S-C Hung. Ontology-based context representation and reasoning using owl and swrl. In *2010 8th Annual Communication Networks and Services Research Conference*, pages 215–220. IEEE, 2010.

Mohan Pammi, Nima Aghaeepour, and Josef Neu. Multiomics, artificial intelligence, and precision medicine in perinatology. *Pediatric Research*, 93:308–315, 2023. doi: 10.1038/s41390-022-02181-x. URL https://doi.org/10.1038/s41390-022-02181-x.

Jiantao Wu, Fabrizio Orlandi, Declan O'Sullivan, and Soumyabrata Dev. Linkclimate: An interoperable knowledge graph platform for climate data. *Computers & Geosciences*, 169:105215, 2022.

Paul Pajo. Leveraging ai-driven hypothesis generation for niche and obscure fields: Applications in musicology. https://www.researchgate.net/publication/390066146_Leveraging_AI-Driven_Hypothesis_Generation_for_Niche_and_Obscure_Fields_Applications_in_Musicology, March 2025. Preprint uploaded to ResearchGate.

Alexander Lavin, David Krakauer, Hector Zenil, Justin Gottschlich, Tim Mattson, Johann Brehmer, Anima Anandkumar, Sanjay Choudry, Kamil Rocki, Atılım Güneş Baydin, et al. Simulation intelligence: Towards a new generation of scientific methods. *arXiv preprint arXiv:2112.03235*, 2021.

Vasanth Sarathy and Matthias Scheutz. Biplex: Creative problem-solving by planning for experimentation. In *International Conference on Computational Creativity*, 2022.

Caron AC Clark, Tomáš Helikar, and Joseph Dauer. Simulating a computational biological model, rather than reading, elicits changes in brain activity during biological reasoning. *CBE—Life Sciences Education*, 19(3):ar45, 2020.

Sébastien Hélie and Ron Sun. Knowledge integration in creative problem solving. In *Proceedings of the 30th Annual Meeting of the Cognitive Science Society*, pages 1681–1686. Austin, TX: Cognitive Science Society, 2008.

Gechun Lin and Christopher Lucas. An introduction to neural networks for the social sciences. 2023.

Jeremy Utley and Perry Klebahn. *Ideaflow: the only business metric that matters*. Penguin, 2022.

Luke de Oliveira, Michela Paganini, and Benjamin Nachman. Learning particle physics by example: location-aware generative adversarial networks for physics synthesis. *Computing and Software for Big Science*, 1(1):4, 2017.

Elahe Khatibi, Mahyar Abbasian, Zhongqi Yang, Iman Azimi, and Amir M Rahmani. Alcm: Autonomous llm-augmented causal discovery framework. *arXiv preprint arXiv:2405.01744*, 2024.

Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of causal inference: foundations and learning algorithms*. The MIT Press, 2017.

Peter JF Lucas. Biomedical applications of bayesian networks. *Advances in probabilistic graphical models*, pages 333–358, 2007.

Leland Gerson Neuberg. Causality: models, reasoning, and inference, by judea pearl, cambridge university press, 2000. *Econometric Theory*, 19(4):675–685, 2003.

Karamarie Fecho, Chris Bizon, Frederick Miller, Shepherd Schurman, Charles Schmitt, William Xue, Kenneth Morton, Patrick Wang, Alexander Tropsha, et al. A biomedical knowledge graph system to propose mechanistic hypotheses for real-world environmental health observations: cohort study and informatics application. *JMIR medical informatics*, 9(7):e26714, 2021.

Yuchen Yan and Chong Chen. Scigraph: A knowledge graph constructed by function and topic annotation of scientific papers.

Sebastian Mosbach, Angiras Menon, Feroz Farazi, Nenad Krdzavac, Xiaochi Zhou, Jethro Akroyd, and Markus Kraft. Multiscale cross-domain thermochemical knowledge-graph. *Journal of Chemical Information and Modeling*, 60(12):6155–6166, 2020.

Emanuele Rossi, Ben Chamberlain, Fabrizio Frasca, Davide Eynard, Federico Monti, and Michael Bronstein. Temporal graph networks for deep learning on dynamic graphs. *arXiv preprint arXiv:2006.10637*, 2020.

Xuemei Gu and Mario Krenn. Forecasting high-impact research topics via machine learning on evolving knowledge graphs, 2025. URL `https://arxiv.org/abs/2402.08640`.

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

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

Alexei Kurakin and Dale E Bredesen. An unconventional iap-binding motif revealed by target-assisted iterative screening (tais) of the bir3-ciap1 domain. *Journal of Molecular Recognition: An Interdisciplinary Journal*, 20(1):39–50, 2007.

Nathaniel H Park, Tiffany J Callahan, James L Hedrick, Tim Erdmann, and Sara Capponi. Leveraging chemistry foundation models to facilitate structure focused retrieval augmented generation in multi-agent workflows for catalyst and materials design. *arXiv preprint arXiv:2408.11793*, 2024.

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

LabKey. Labkey server, 2025. URL `https://www.labkey.org/`.

MathWorks. *Simulink: Simulation and Model-Based Design*. The MathWorks, Inc., Natick, MA, 2025. URL `https://www.mathworks.com/products/simulink.html`.

F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. scikit-learn: Machine learning in python, 2011. URL `https://scikit-learn.org/`.

Sean J Taylor and Benjamin Letham. Forecasting at scale. *The American Statistician*, 72(1):37–45, 2018.

TensorFlow. Tensorflow serving, 2025. URL `https://www.tensorflow.org/tfx/guide/serving`.

IBM. Ibm watson, 2025. URL `https://www.ibm.com/watson`.

The Zooniverse Team. Zooniverse: People-powered research, 2025. URL `https://www.zooniverse.org/`.

Project Jupyter. Project jupyter: Jupyter notebooks, 2025. URL `https://jupyter.org/`.

Center for Causal Discovery. Tetrad: Causal discovery software, 2025. URL `https://www.ccd.pitt.edu/tools/tetrad/`.

MLPerf. Mlperf: Fair and useful benchmarks for machine learning, 2025. URL `https://mlperf.org/`.

Yoshua Bengio, Tristan Deleu, Nasim Rahaman, Rosemary Ke, Sébastien Lachapelle, Olexa Bilaniuk, Anirudh Goyal, and Christopher Pal. A meta-transfer objective for learning to disentangle causal mechanisms. *arXiv preprint arXiv:1901.10912*, 2019.

Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and S Yu Philip. A comprehensive survey on graph neural networks. *IEEE transactions on neural networks and learning systems*, 32(1):4–24, 2020.

Judea Pearl. Causal inference in statistics: An overview. 2009.

Paul R Rosenbaum and Donald B Rubin. The central role of the propensity score in observational studies for causal effects. *Biometrika*, 70(1):41–55, 1983.

Joshua D Angrist, Guido W Imbens, and Donald B Rubin. Identification of causal effects using instrumental variables. *Journal of the American statistical Association*, 91(434):444–455, 1996.

Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. *Advances in neural information processing systems*, 29, 2016.

Paul Pu Liang, Irene Mengze Li, Emily Zheng, Yao Chong Lim, Ruslan Salakhutdinov, and Louis-Philippe Morency. Towards debiasing sentence representations. *arXiv preprint arXiv:2007.08100*, 2020.

Jwala Dhamala, Tony Sun, Varun Kumar, Satyapriya Krishna, Yada Pruksachatkun, Kai-Wei Chang, and Rahul Gupta. Bold: Dataset and metrics for measuring biases in open-ended language generation. In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*, pages 862–872, 2021.