

AI4Research: A Survey of Artificial Intelligence for Scientific Research

Qiguang Chen^{1*} Mingda Yang^{1*} Libo Qin^{2,✉} Jinhao Liu¹ Zheng Yan¹ Jiannan Guan¹
 Dengyun Peng¹ Yiyang Ji¹ Hanjing Li¹ Mengkang Hu³ Yimeng Zhang⁴ Yihao Liang⁴
 Yuhang Zhou⁵ Jiaqi Wang⁶ Zhi Chen⁷ Wanxiang Che^{1,✉}

¹ LARG, Research Center for Social Computing and Interactive Robotics, Harbin Institute of Technology,

² School of Computer Science and Engineering, Central South University, ³ The University of Hong Kong,

⁴ Independent Researcher, ⁵ Fudan University, ⁶ Chinese University of Hong Kong,

⁷ ByteDance Seed (China)

Abstract:

Recent advancements in artificial intelligence (AI), particularly in large language models (LLMs) such as OpenAI-o1 and DeepSeek-R1, have demonstrated remarkable capabilities in complex domains such as logical reasoning and experimental coding. Motivated by these advancements, numerous studies have explored the application of AI in the innovation process, particularly in the context of scientific research. These AI technologies primarily aim to develop systems that can autonomously conduct research processes across a wide range of scientific disciplines. Despite these significant strides, a comprehensive survey on AI for Research (AI4Research) remains absent, which hampers our understanding and impedes further development in this field. To address this gap, we present a comprehensive survey and offer a unified perspective on AI4Research. Specifically, the main contributions of our work are as follows: (1) **Systematic taxonomy**: We first introduce a systematic taxonomy to classify five mainstream tasks in AI4Research. (2) **New frontiers**: Then, we identify key research gaps and highlight promising future directions, focusing on the rigor and scalability of automated experiments, as well as the societal impact. (3) **Abundant applications and resources**: Finally, we compile a wealth of resources, including relevant multidisciplinary applications, data corpora, and tools. We hope our work will provide the research community with quick access to these resources and stimulate innovative breakthroughs in AI4Research.

* *Equal Contribution*

✉ *Corresponding Author*

Keywords: AI4Research, Large Language Models, Scientific Comprehension, Academic Survey, Scientific Discovery, Academic Writing, Academic Peer Review

 **Date:** July 02, 2025

 **Projects:** <https://ai-4-research.github.io>

 **Code Repository:** <https://github.com/LightChen233/Awesome-AI4Research>

 **Contact:** qgchen@ir.hit.edu.cn, car@ir.hit.edu.cn, lbqin@csu.edu.cn

Contents

1	Introduction	5
2	The Definition of AI4Research	6
2.1	Component-wise Definition of AI4Research	8
2.1.1	AI for Scientific Comprehension	8
2.1.2	AI for Academic Survey	8
2.1.3	AI for Scientific Discovery	9
2.1.4	AI for Academic Writing	9
2.1.5	AI for Academic Peer Reviewing	9
2.2	Discussion About AI4Science and AI4Research	10
3	AI for Scientific Comprehension	11
3.1	Textual Scientific Comprehension	11
3.1.1	Semi-Automatic Scientific Comprehension	11
3.1.2	Full-Automatic Scientific Comprehension	12
3.2	Table & Chart Scientific Comprehension	13
3.2.1	Table Understanding	13
3.2.2	Chart Understanding	13
4	AI for Academic Survey	13
4.1	Related Work Retrieval	14
4.2	Overview Report Generation	15
4.2.1	Research Roadmap Mapping	15
4.2.2	Section-level Related Work Generation	15
4.2.3	Document-level Survey Generation	16
5	AI for Scientific Discovery	17
5.1	Idea Mining	17
5.1.1	Idea Mining from Internal Knowledge	17
5.1.2	Idea Mining from External Signal	18
5.1.3	Idea Mining from Team discussion	19

5.2	Novelty & Significance Assessment	20
5.3	Theory Analysis	20
5.3.1	Scientific Claim Formalization	20
5.3.2	Scientific Evidence Collection	21
5.3.3	Scientific Verification Analysis	21
5.3.4	Theorem Proving	21
5.4	Scientific Experiment Conduction	21
5.4.1	Experiment Design	22
5.4.2	Pre-Experiment Estimation	22
5.4.3	Experiment Management	23
5.4.4	Experimental Conduction	23
5.4.5	Experimental Analysis	25
5.5	Full-Automatic Discovery	25
6	AI for Academic Writing	25
6.1	Semi-Automatic Academic Writing	26
6.1.1	Assistance During Manuscript Preparation	26
6.1.2	Assistance During Manuscript Writing	26
6.1.3	Assistance After Manuscript Completion	27
6.2	Full-Automatic Academic Writing	28
7	AI for Academic Peer Reviewing	28
7.1	Pre-Review	28
7.1.1	Desk-Review	29
7.1.2	Reviewer Matching	29
7.2	In-Review	29
7.2.1	Peer-Review	30
7.2.2	Meta-Review	31
7.3	Post-Review	31
7.3.1	Influence Analysis	31
7.3.2	Promotion Enhancement	32
8	Application of AI for Research	32

8.1	AI for Natural Science Research	33
8.1.1	AI for Physics Research	33
8.1.2	AI for Biology & Medical Research	33
8.1.3	AI for Chemistry & Materials Research	35
8.2	AI for Applied Science and Engineering Research	36
8.2.1	AI for Robotics and Control Research	36
8.2.2	AI for Software Engineering	37
8.3	AI for Social Science Research	37
8.3.1	AI for Sociology Research	37
8.3.2	AI for Psychology Research	38
9	Resources	39
9.1	AI for Scientific Comprehension	39
9.1.1	Textual Scientific Comprehension	39
9.1.2	Table & Chart Scientific Comprehension	40
9.2	AI for Academic Survey	40
9.3	AI for Scientific Discovery	40
9.4	AI for Academic Writing	43
9.4.1	Semi-Automatic Academic Writing	43
9.5	AI for Academic Peer Reviewing	43
10	Frontiers & Future Direction	44
10.1	Interdisciplinary AI Models	44
10.2	Ethics and Safety in AI4Research	45
10.3	AI for Collaborative Research	45
10.4	Explainability and Transparency of AI4Research	46
10.5	AI for Dynamic and Real-Time Optimized Scientific Experimentation	46
10.6	Multimodal Integration in AI4Research	47
10.7	Multilingual Integration in AI4Research	47
11	Related work	48
12	Conclusion	48

1. Introduction

In recent years, the rise of artificial intelligence (AI), particularly large language models (LLMs) like DeepSeek-R1 [263], has stimulated significant research in the field of reasoning. These breakthroughs have notably improved the models' performance across diverse areas, including mathematical reasoning, programming, and interdisciplinary knowledge [724, 748, 616, 931, 947, 110]. Some of these models have even surpassed the Turing Test [352], marking a pivotal achievement in AI development. Inspired by these, a series of works attempts to explore advanced AI systems for innovative tasks, especially in the scientific discovery of new research [863, 887, 847, 948]. Earlier, the AI Scientist [507] introduces the concept of a fully automatic AI for research system, which divides the research process into three key stages: idea mining, experiment conduction, and academic writing. Initially, the system generates and evaluates novel ideas and hypotheses. Once a hypothesis is formulated, experiments are conducted automatically, producing results that include numerical data and visual summaries. These results are presented in tables and images, followed by an interpretation with a convincing description, culminating in a LaTeX report. In the final stage, the AI Scientist generates an automated review that refines the project and provides feedback for future open scientific discoveries. Similarly, other classic models, such as Carl [330] and Zochi [12], follow broadly analogous workflows. Notably, AgentArxiv [665] and AgentLab [666] assign distinct roles to multiple agents to simulate the collaborative nature of scientific research teams, incorporating additional peer review, academic survey, enabling semi-automatic and even full-automatic collaboration rather than relying on a single agent [478, 870, 112, 658, 53]. Despite these advancements, there remains a lack of comprehensive

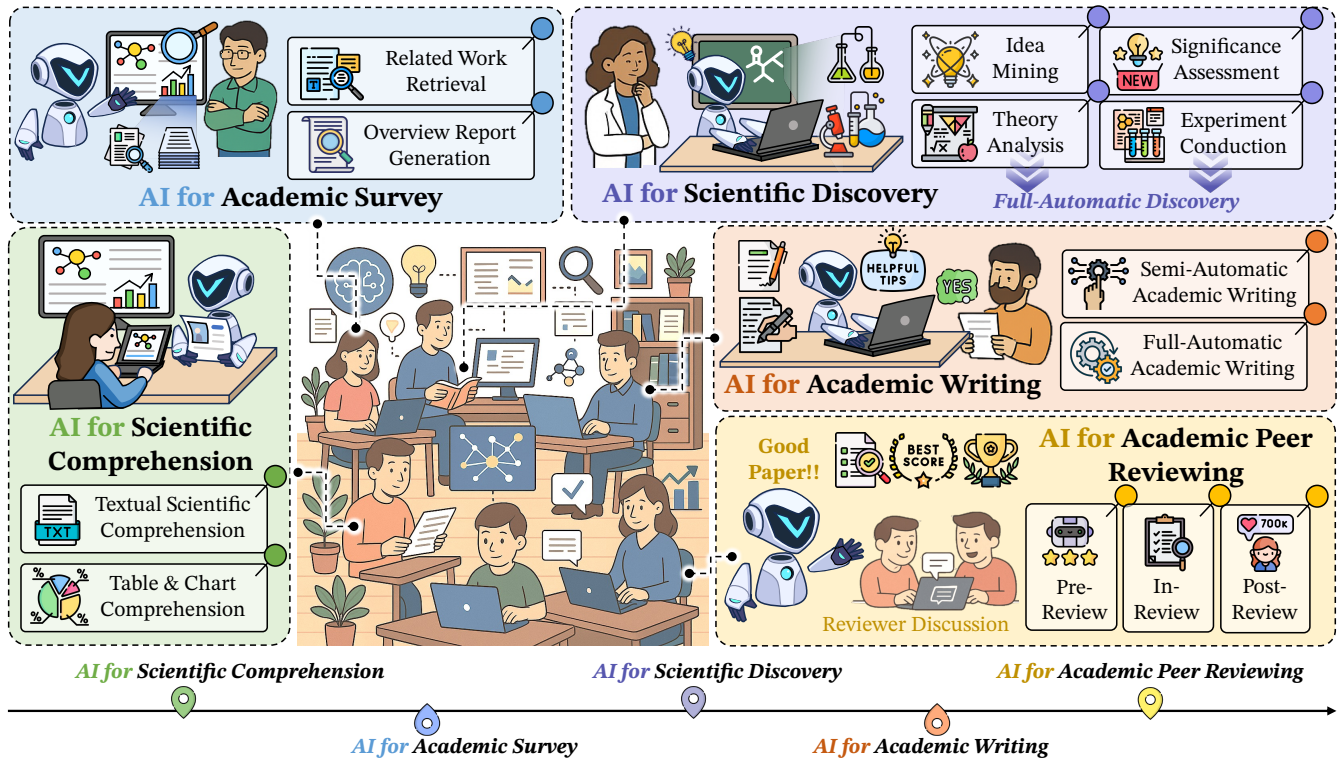


Figure 1: The mainstream processes and categories of AI4Research, which can be divided into five key areas: (1) AI for Scientific Comprehension, (2) AI for Academic Survey, (3) AI for Scientific Discovery, (4) AI for Academic Writing, and (5) AI for Academic Peer Review. Each of these areas contributes to improving the effectiveness and efficiency of AI-integrated research and publication.

surveys to systematically analyze the key factors and recent developments in AI-driven research, which significantly impedes the continued progress of this field.

To address this gap, we first define and present a comprehensive survey of AI for research, termed AI4Research. As shown in Figure 1, we **introduce a systematic taxonomy of AI4Research**, focusing on the following areas: (1) *AI for Scientific Comprehension*: AI systems’ ability to extract relevant information from scientific literature is crucial; (2) *AI for Academic Surveys*: This involves AI techniques for systematically reviewing and summarizing scientific literature; (3) *AI for Scientific Discovery*: AI is used to generate hypotheses, theories, or models based on existing scientific knowledge; (4) *AI for Academic Writing*: AI tools support researchers in drafting, editing, and formatting manuscripts; (5) *AI for Academic Reviewing*: AI assists in evaluating and providing feedback on scientific manuscripts. Given the vast literature, we **highlight promising future research in AI4Research**. Future work should prioritize interdisciplinary AI models that integrate knowledge across scientific domains to encourage cross-disciplinary collaboration. Addressing ethical concerns and biases within AI systems is crucial for ensuring fairness and transparency in research. Improving the explainability of AI models and exploring adaptive, real-time systems for dynamic scientific experiments will be vital for advancing AI’s role in research. Additionally, we **suggest key applications and valuable resources in AI4Research**, such as representative multidisciplinary applications, open-source frameworks, and datasets repositories to support further studies. We introduce AI for Natural Science research, AI for Applied Science and Engineering research, and AI for Social Science research. Finally, we review tools essential for model development and public benchmarks that provide rich data for training and experimentation.

The main contributions of this work are as follows:

- **Systematic Taxonomy for AI in Research**: This paper introduces a comprehensive taxonomy of AI applications in research, spanning five areas: scientific comprehension, academic surveys, scientific discovery, academic writing, and academic reviewing. It categorizes AI tools that enhance and even automatically execute various stages of the research process.
- **Emerging Future Research Areas**: The paper identifies key future research avenues for AI in academia, including the development of interdisciplinary AI models, addressing ethical concerns and biases, improving model explainability, and exploring adaptive AI systems for dynamic scientific experiments.
- **Key Applications and Abundant Trending Resources**: We present multidisciplinary AI4Research applications across natural sciences, applied science, and social sciences. It also identifies essential resources, open-source frameworks, public datasets, collaborative platforms, cloud-based AI services, and academic tools, that facilitate discovery management, data processing, and AI-driven research.

2. The Definition of AI4Research

AI4Research denotes the application of artificial intelligence methods to improve, accelerate, and partially automate research across disciplines. To clarify this paradigm, as shown in Figure 2, we identify six core capabilities: AI for Scientific Comprehension, AI for Academic Survey, AI for Scientific Discovery, AI for Academic Writing, AI for Academic Peer Reviewing. Each of them illustrates a distinct way that AI advances the research process. Formally, let $\mathcal{T} = \{T_{SC}, T_{AS}, T_{SD}, T_{AW}, T_{PR}\}$ be the set of research tasks, Scientific Comprehension, Academic Survey, Scientific Discovery, Academic Writing, and Peer Reviewing. For each task $T_i \in \mathcal{T}$, there exists a corresponding AI model A_i that is specifically tailored to address the requirements of that task. Then the overall AI4Research system can be expressed as the functional composition:

$$\mathcal{A} = A_{PR} \circ A_{AW} \circ A_{SD} \circ A_{AS} \circ A_{SC}, \quad (1)$$

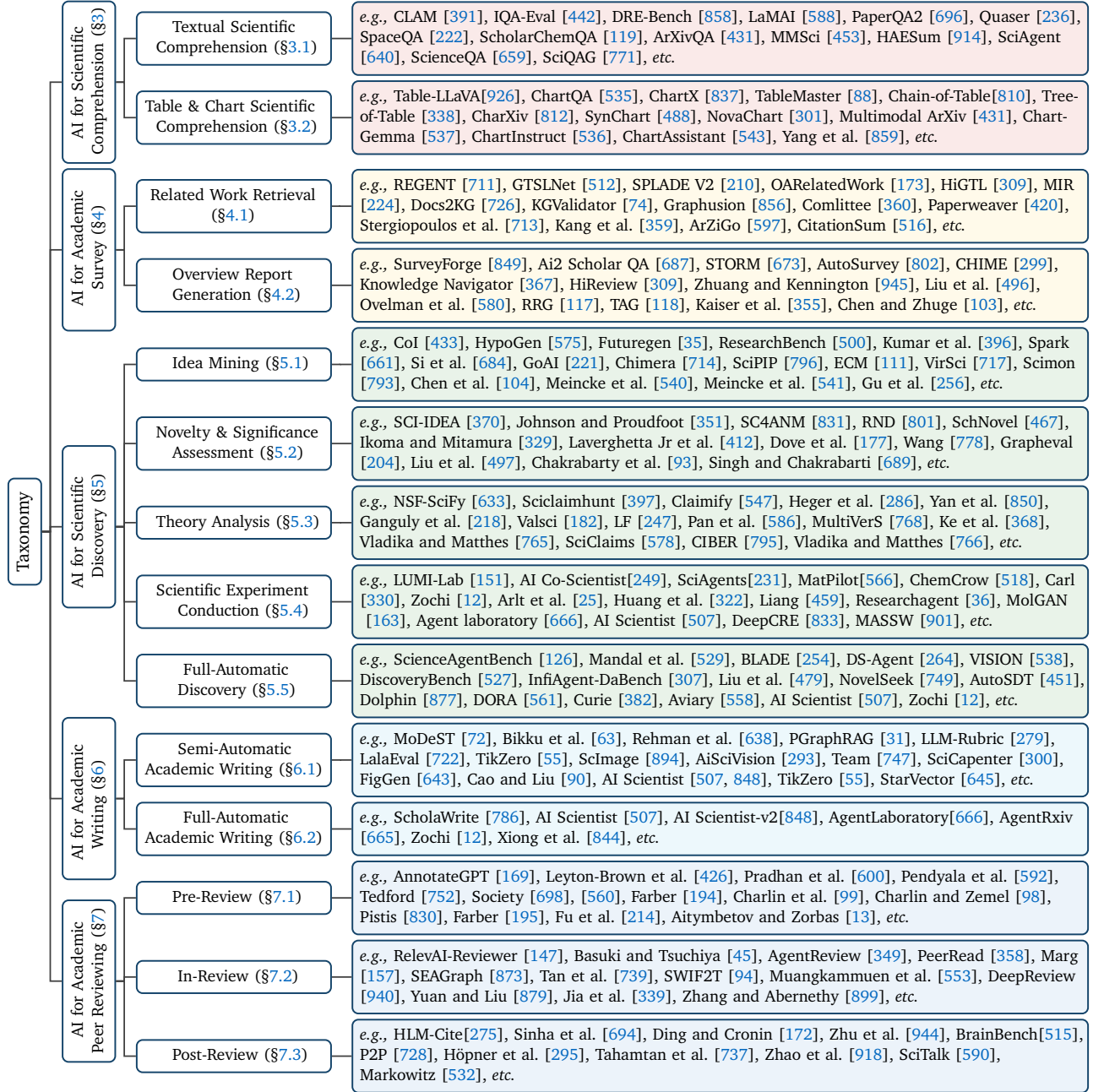


Figure 2: The taxonomy of AI in research (AI4Research) is categorized into five key areas. Each area is subdivided into specific tasks, underscoring the varied roles of AI in the entire research process.

where \circ denotes the function composition operator, meaning that the output of one function becomes the input of the next. Further, applied to a research query q (or more generally to an interactive research-query lifecycle Q), we can obtain:

$$\mathcal{A}(q) = (A_{PR} \circ A_{AW} \circ A_{SD} \circ A_{AS} \circ A_{SC})(q). \quad (2)$$

The objective of an AI4Research system is to maximize research lifecycle efficiency, application perfor-

mance, and innovation capacity, namely:

$$\max \{ \eta(\mathcal{A}(\mathcal{Q})), \alpha(\mathcal{A}(\mathcal{Q})), \tau(\mathcal{A}(\mathcal{Q})) \}, \quad (3)$$

where $\eta(\cdot)$, $\tau(\cdot)$, and $\alpha(\cdot)$ evaluates the efficiency, performance, and innovation of the generated research publications $\mathcal{A}(\mathcal{Q})$, respectively.

2.1. Component-wise Definition of AI4Research

We now define and formalize each core module in the AI4Research framework.

2.1.1. AI for Scientific Comprehension

AI for Scientific Comprehension (AI4SC) is central to AI4Research, enabling extraction, interpretation, and synthesis of information from a single scientific literature. This accelerates human knowledge acquisition and improves the efficiency of automated analysis. Formally, we define this module to gain knowledge K after comprehension as a composite reasoning function:

$$\hat{\mathcal{K}} = A_{SC}(\mathcal{K}) = f_{SC}(\mathcal{K}|D_{SC}, \Phi_{SC}) = f_{TCSC} \circ f_{TSC}(\mathcal{K}|D_{SC}, \Phi_{SC}), \quad (4)$$

where A_{SC} is the comprehension AI model to extract the possible knowledge \mathcal{K} ; the documents $D_{SC} = \{D_T, D_F, D_M\}$ comprises texts (D_T), figures (D_F), and other metadata (D_M); Φ_{SC} includes model parameters and domain priors; and f_{SC} means the specific comprehension algorithms, including a textual comprehension function f_{TSC} that extracts and interprets textual content, and a table & chart comprehension function f_{TCSC} that analyzes tables and charts.

The goal of AI4SC is to maximize scientific understanding σ with extracted knowledge $\hat{\mathcal{K}}$ from the original documents D_{SC} :

$$\max\{\sigma\} = \max\{\mathbb{E}_{\hat{\mathcal{K}} \sim A_{SC}}[\text{Coherence}(\hat{\mathcal{K}}, D_{SC}) + \text{Coverage}(\hat{\mathcal{K}}, D_{SC})]\}, \quad (5)$$

where Coherence measures logical consistency; Coverage quantifies concept completeness between them.

2.1.2. AI for Academic Survey

AI for Academic Survey (AI4AS) is designed to synthesize and structure multiple existing literature, providing a comprehensive overview of a research domain. This enhances the ability to identify trends, gaps, and key contributions in scientific fields. Formally, we define this module to generate a structured literature survey S as a functional synthesis function:

$$\hat{\mathcal{S}} = A_{AS}(\mathcal{S}) = f_{AS}(\mathcal{S}|R_{AS}, \Phi_{AS}) = f_{Gen} \circ f_{Retrieval}(\mathcal{S}|R_{AS}, \Phi_{AS}), \quad (6)$$

where A_{AS} is the survey AI model to generate the possible survey \mathcal{S} ; R_{AS} comprising survey domain requirements; Φ_{AS} includes model parameters and domain priors; f_{AS} means the specific survey algorithms, which include a retrieval function $f_{Retrieval}$ that retrieves relevant literature based on the query, and a generative function f_{Gen} that produces thematic clusters and summaries.

The objective of AI4AS is to maximize survey quality ρ of the generated survey $\hat{\mathcal{S}}$ with respect to the requirement R_{AS} :

$$\max\{\rho\} = \max\{\mathbb{E}_{\hat{\mathcal{S}} \sim A_{AS}}[\text{Relevance}(\hat{\mathcal{S}}, R_{AS}) + \text{Coverage}(\hat{\mathcal{S}}, R_{AS}) + \text{Clarity}(\hat{\mathcal{S}}, R_{AS})]\}, \quad (7)$$

where Relevance measures the match between documents and the target topic; Coverage assesses the breadth and depth of the domain; Clarity reflects the coherence, abstraction quality, and utility of the synthesized representation based on the generated survey and requirements.

2.1.3. AI for Scientific Discovery

AI for Scientific Discovery (AI4SD) is focused on generating, and validating novel scientific hypotheses or ideas and conducting experiments or simulations. This module enhances the ability to explore uncharted scientific territories and accelerate innovation. Formally, we define this module to generate, validate, and implement scientific innovations $\hat{\mathcal{I}}$ as a discovery-oriented function:

$$\hat{\mathcal{I}} = A_{SD}(\mathcal{I}) = f_{SD}(\mathcal{I}|K_{SD}, R_{SD}, \Phi_{SD}) = f_{ED} \circ f_{TA} \circ f_{NSA} \circ f_{IM}(\mathcal{I}|K_{SD}, R_{SD}, \Phi_{SD}), \quad (8)$$

where A_{SD} is a discovery-oriented AI to explore possible innovation \mathcal{I} ; scientific knowledge $K_{SD} = \{K_D, K_{AS}\}$ is the given domain knowledge (K_D) and recent related-work summarized knowledge (K_{AS}) derived from upstream comprehension and survey stages; R_{SD} means the research requirement; Φ_{SD} includes model parameters and domain priors; f_{SD} means the specific discovery algorithms, which include a generative function f_{IM} that mines candidate ideas, a novelty and significance assessment function f_{NSA} that evaluates the quality and importance of each idea candidates, a theory analysis function f_{TA} that checks theoretical soundness, and an experiment conduction function f_{ED} that makes plans and executes experiments then finally complete the scientific discovery.

The goal of AI4SD is to maximize the total discovery quality δ of the generated innovations $\hat{\mathcal{I}}$:

$$\max\{\delta\} = \max\{\mathbb{E}_{\hat{\mathcal{I}} \sim A_{SD}}[\text{Novelty}(\hat{\mathcal{I}}) + \text{Validity}(\hat{\mathcal{I}}) + \text{Significance}(\hat{\mathcal{I}})]\}, \quad (9)$$

where Novelty evaluates innovativeness; Validity assesses experimental and theoretical soundness; Significance reflects the follow-up impact of the study.

2.1.4. AI for Academic Writing

AI for Academic Writing (AI4AW) is a highlight section of AI4Research, assisting researchers in generating, revising, and formatting scientific manuscripts. This module enhances the quality and efficiency of academic writing, ensuring that manuscripts are well-structured and compliant with publication standards. Formally, we define this module to generate a publication-ready manuscript \mathcal{M} as a collaborative writing function:

$$\hat{\mathcal{M}} = A_{AW}(\mathcal{M}) = f_{AW}(\mathcal{M}|K_{AS}, \text{Info}_I, \Phi_{AW}) = f_{DWP} \circ f_{DMW} \circ f_{AWC}(\mathcal{M}|K_{AS}, \text{Info}_I, \Phi_{AW}), \quad (10)$$

where A_{AW} denotes a writing-oriented AI to generate the possible manuscript \mathcal{M} ; Info_I is all information in the scientific discovery stage, including ideas, experimental designs, and attachments such as codes and data; Φ_{AW} includes model parameters and domain priors; f_{AW} means the specific writing algorithms, which include a during-manuscript-preparation function f_{DWP} that prepares the manuscript structure, a during-manuscript-writing function f_{DMW} that generates the manuscript content, and a after-manuscript-completion function f_{AWC} that completes grammatical corrections, expressions and logical modifications.

The objective of AI4AW is to maximize writing quality and effectiveness ω of the manuscript $\hat{\mathcal{M}}$:

$$\max\{\omega\} = \max\{\mathbb{E}_{\hat{\mathcal{M}} \sim A_{AW}}[\text{Consistency}(\hat{\mathcal{M}}) + \text{Readability}(\hat{\mathcal{M}}) + \text{Compliance}(\hat{\mathcal{M}})]\}, \quad (11)$$

where Consistency reflects logical flow and internal coherence; Readability measures linguistic clarity and ease of understanding; Compliance assesses adherence to formatting and stylistic requirements.

2.1.5. AI for Academic Peer Reviewing

AI for Academic Peer Reviewing (AI4PR) is a critical component of AI4Research, automating and enhancing the peer review process. This module aims to provide structured, objective, and constructive reviews of







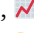


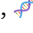





	AI4Science	AI4Research
Scope	 Scientific Discovery,  Data Analysis.	 Broader Research workflows.
Goal	 Scientific Breakthroughs.	 Publications,  Methods,  Overall Productivity.
Applications	 Material Discovery,  Drug Design,  Genomics, <i>etc.</i>	 Comprehension, Writing,  Peer Review, <i>etc.</i>
Target Users	 Research Experts.	Both  Research Experts and  New Scientists.

Table 1: Comparison and discussion between AI4Science and AI4Research, especially in terms of scope, goal, applications, and target users.

scientific manuscripts, improving the quality and efficiency of the review cycle. Formally, we define this module to generate a structured review result R as an evaluative reasoning function:

$$\hat{\mathcal{R}} = A_{PR}(\mathcal{R}) = f_{PR}(\mathcal{R}|P, \Phi_{PR}) = f_{PostP} \circ f_{InP} \circ f_{PreP}(\mathcal{R}|\hat{\mathcal{M}}, \Phi_{PR}), \quad (12)$$

where A_{PR} denotes a review-oriented AI to generate the possible review \mathcal{R} ; Φ_{PR} includes model parameters and domain priors; f_{PR} means the specific review algorithms, which include a pre-review function f_{PreP} that completes pre-review preparations, an in-review function f_{InP} that generates or augments review reports, and a post-review function f_{PostP} that completes post-review analysis of papers.

The goal of AI4PR is to maximize review quality θ of the review result $\hat{\mathcal{R}}$ based on the manuscript $\hat{\mathcal{M}}$:

$$\max\{\theta\} = \max\{\mathbb{E}_{\hat{\mathcal{R}} \sim A_{PR}}[\text{Correctness}(\hat{\mathcal{R}}, \hat{\mathcal{M}}) + \text{Helpfulness}(\hat{\mathcal{R}}, \hat{\mathcal{M}}) + \text{Consistency}(\hat{\mathcal{R}}, \hat{\mathcal{M}})]\}, \quad (13)$$

where Correctness means the review can correctly reflect the pros and cons of research; Helpfulness measures the depth, constructiveness, and usefulness of feedback; Consistency quantifies the alignment of the review with established evaluation criteria and domain standards.

2.2. Discussion About AI4Science and AI4Research

Since the concepts of AI4Science and AI4Research share many similarities, as outlined in Table 1, it is important to distinguish the key differences between the two. **AI4Science (Artificial Intelligence for Science)** focuses on applying AI technologies to accelerate scientific discovery and data analysis across various fields, including material discovery, drug design, and genomic analysis. Its primary objective is to integrate AI into research workflows to support experts in achieving significant scientific advancements. In contrast, **AI4Research (Artificial Intelligence for Research)** adopts a broader perspective, addressing publications, methodologies, and overall research productivity. It emphasizes AI’s role in enhancing research methods and supporting the academic environment for both established researchers and emerging scientists. Key applications in this domain include AI-driven tools for literature comprehension, academic writing assistance, and peer review processes.

The core distinction between these frameworks lies in their focus: AI4Science targets specific scientific problems and experimental protocols, while AI4Research addresses broader research methodologies and academic infrastructure. However, as LLMs develop more advanced reasoning and generative capabilities, a unified workflow is emerging that can address both specialized scientific challenges and general research processes. Consequently, AI4Science tools are increasingly integrated into AI4Research environments, often serving as callable components in LLM-based systems for scientific exploration. Subsequently, we will provide a detailed analysis of our taxonomy and the relevant literature.

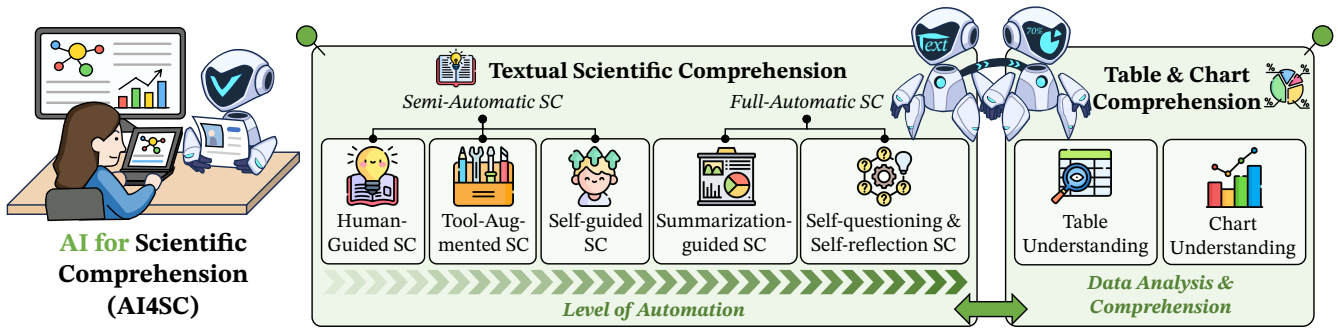


Figure 3: The primary paradigms of AI for Scientific Comprehension. These include: (1) Textual Scientific Comprehension, which is further categorized into Semi-Automatic and Fully-Automatic Scientific Comprehension; and (2) Table & Chart Scientific Comprehension, encompassing Table and Chart Understanding.

3. AI for Scientific Comprehension

Scientific comprehension plays a pivotal role in advancing AI4Research, encompassing the ability to extract, understand, and synthesize information from the scientific literature. This capability not only accelerates human understanding and knowledge acquisition but also enhances the efficiency of automatic analysis, enabling more effective research processing. As shown in Figure 3, it contains two main categories: Textual Scientific Comprehension (§ 3.1) and Table & Chart Scientific Comprehension (§ 3.2).

3.1. Textual Scientific Comprehension

Textual Scientific Comprehension refers to the ability to understand, interpret, and critically evaluate scientific texts. It involves identifying key concepts, grasping complex terminology, and synthesizing information to form a cohesive understanding of scientific principles and findings [621, 70, 531]. As depicted in Figure 3 (middle), we categorize the corresponding comprehension technologies into two types based on automation level: Semi-Automatic and Fully-Automatic Scientific Comprehension.

3.1.1. Semi-Automatic Scientific Comprehension

Semi-automatic scientific comprehension denotes systems in which, given a manually created question, the AI produces comprehensive question-related comprehension of long-context scientific content. Such systems support both researchers and AI models in deepening their grasp of complex scientific concepts [119, 290, 818]. Specifically, these systems comprise three main categories:

Human-Guided Scientific Comprehension is an interactive approach where researchers and language models engage in iterative dialogues to produce a deepened understanding of questions on complex scientific literature step-by-step [391, 895, 442, 858, 613]. LaMAI [588] equips language models with “active inquiry” capabilities: before providing a definitive answer, the model asks clarifying questions to resolve ambiguities in user queries, reducing misinterpretations and enhancing relevance [848]. These platforms illustrate that embedding structured human feedback loops within LLM-based tools improves output reliability and enriches the scientific discovery process by uncovering latent questions and assumptions. However, the approach requires significant human-AI interaction, which can increase costs.

Tool-Augmented Scientific Comprehension refers to cases where a researcher’s query surpasses a language model’s knowledge base or its context-window limit [822]. The model then invokes several external tools to ensure accurate output: (1) *Knowledge Retrieval Tool* uses retrieval-augmented generation to inject

knowledge beyond the model’s training [378, 659]. Early systems like document-centric agents [405] extract key findings, note limitations, and propose future directions. Graphusion [856] advances this with a zero-shot RAG approach: it builds scientific knowledge graphs by extracting entity triples, merging duplicates, and resolving conflicts across disciplines—without manual effort. SiGIR [135] uses self-critique feedback to guide the iterative reasoning process during knowledge-intensive multi-hop reasoning tasks. **(2) Fact Checking Tool** mitigates hallucinations and factual errors by applying verification modules to reduce the AI’s hallucinations [354, 278, 248, 909]. PaperQA2 [696] integrates rigorous factuality checks and matches or exceeds expert accuracy on literature-review tasks, all without unrestricted Internet access or human oversight. **(3) Reasoning-Augmentation Tool** addresses limited logical reasoning and computation in standalone models to deepen the AI’s theory-level comprehension [110]. For example, SciAgent [522] dynamically selects calculators and formula evaluators to deliver precise, domain-specific reasoning. Collectively, these advances show how coupling language models with specialized tools transforms scientific workflows from passive consumption into an interactive, tool-powered process that accelerates discovery while preserving rigor.

Self-guided Scientific Comprehension refers to a model’s capacity to respond to a single-turn query regarding a scientific publication with a comprehensive, context-sensitive answer [56, 650, 71]. Earlier, Clark et al. [142] demonstrate that even seemingly factual questions about academic papers require deep contextual understanding and meticulous attention to document-specific details [236]. To address these challenges, subsequent studies focus more on enhanced semi-automatic scientific comprehension in long-context papers [487, 743], particularly in specialized fields such as aerospace science [222], chemistry [594, 119], and clinical medicine [636, 693]. It illustrates that enhancing models to align with the linguistic and conceptual conventions of each discipline, particularly those with improved long-context capabilities, leads to significant advancements [125, 498]. Furthermore, recognizing the inherently multimodal nature of scientific papers, several studies have begun to integrate textual analysis with figures and charts [431, 453, 640], thereby advancing towards a more holistic, paper-wide comprehension of scientific content.

3.1.2. Full-Automatic Scientific Comprehension

AI for full-automatic scientific comprehension refers to the ability of an AI system to read and understand scientific knowledge independently without human questions or other intervention. The goal of such systems is to fully automate the processing of scientific literature, the formulation and answering of complex questions, and even, to some extent, scientific discovery or idea mining.

Summarization-guided Automatic Scientific Comprehension refers to the capability of LLMs to autonomously generate summaries of scientific articles and, based on these summaries, construct a comprehensive narrative of the research [369, 209]. This process enhances the model’s holistic understanding of lengthy scientific texts and mitigates comprehension biases that arise from processing extensive documents in a purely token-by-token manner [914]. Furthermore, Ifargan et al. [328] suggest that LLMs can further enhance their overall comprehension of lengthy scientific documents through the generation of autonomous summaries. Their approach utilizes a system of multiple agents, such as a Summary Agent and a Proofreading Agent, working collaboratively to extract and refine experimental results and research methodologies without human intervention. This ultimately produces a refined abstract suitable for peer review.

Self-Questioning & Self-Reflection Automatic Scientific Comprehension involve an AI generating and answering its own questions or reflection to deepen its understanding of scientific content [311, 548, 869]. Earlier, SciInstruct [889] proposes a self-reflective annotation framework, where a model generates step-by-step reasoning for unlabeled scientific questions and then refines its output through self-critique, producing high-quality annotations. Building on this, several studies [514] have focused on prompting models to autonomously create question sequences that enhance their comprehension of scientific texts. One notable

example is SciQAG [771], which proposes a pipeline where a “question generator” and an “answer evaluator” collaborate to extract diverse, research-level comprehension from scientific papers.

More recently, LLMs have been directed to self-improve by posing clarifying questions and decomposing complex problems in a Socratic style, strengthening reasoning and conceptual understanding [622, 705, 811]. The Introspective Growth framework [829] further refines this approach, prompting smaller models to generate fundamental, open-ended questions that guide larger models toward better task comprehension. This process integrates external text retrieval to refine the understanding of technical semantics.

3.2. Table & Chart Scientific Comprehension

Beyond pure textual content, LLMs are employing various techniques to more efficiently interpret and leverage information from tables and figures, thereby achieving a deeper and more comprehensive understanding of scientific literature [140, 562, 266].

3.2.1. Table Understanding

Table understanding involves methods that enable LLMs to extract, interpret, and infer data from tables in scientific literature [718, 832, 28]. **(1) Data Augmentation:** The most direct way is to add higher quality table understanding data. For instance, Zheng et al. [926] introduce the MMTab dataset for large-scale multimodal table understanding in a generative format and propose Table-LLaVA, which reasons directly on table images through instruction tuning, demonstrating the great advantage of visually grounded table representations. **(2) Reasoning Paradigm Augmentation:** Subsequent work explores suitable reasoning paradigms [88, 904, 773, 774]. Wang et al. [810] propose Chain-of-Table, which incrementally constructs and updates tables within an LLM’s reasoning chain to improve comprehension of complex tables. Ji et al. [338] introduce Tree-of-Table, hierarchically condensing and decomposing large tables into a tree structure to facilitate LLM reasoning. Cao and Liu [88] present TableMaster, a framework that enhances LLM table understanding by extracting and verbalizing relevant table content with enriched semantic context and adaptively switching between textual and symbolic reasoning.

3.2.2. Chart Understanding

Chart Understanding involves techniques enabling multimodal large language models to directly process and interpret chart images in scientific papers, supporting tasks such as question answering and summarization based on chart content [602, 463, 544, 319]. Furthermore, several studies focus on assembling and synthesizing large, diverse chart datasets to improve chart understanding [488, 301, 431]. Masry et al. [536] and Meng et al. [543] present vision-language instruction datasets for charts, and train both end-to-end and pipeline models that achieve state-of-the-art results on scientific chart understanding tasks [537]. Further, Yang et al. [859] propose the Formalized Description for Visualization (FDV), a structured textual representation of charts that enables large language models to learn for diverse and deeper comprehension.

4. AI for Academic Survey

It is widely acknowledged that a thorough and well-conducted pre-writing survey and research phase forms the cornerstone of a successful academic article [646]. Inspired by this, AI for Academic Survey is proposed to systematically review and summarize scientific literature through the application of artificial intelligence techniques. This process plays a crucial role in ensuring that researchers and automated systems remain current with the latest advancements in their field and can efficiently identify relevant studies to inform their own work. As shown in Figure 4, it contains two main stages: Related Work Retrieval (§ 4.1) and Overview Report Generation (§ 4.2).

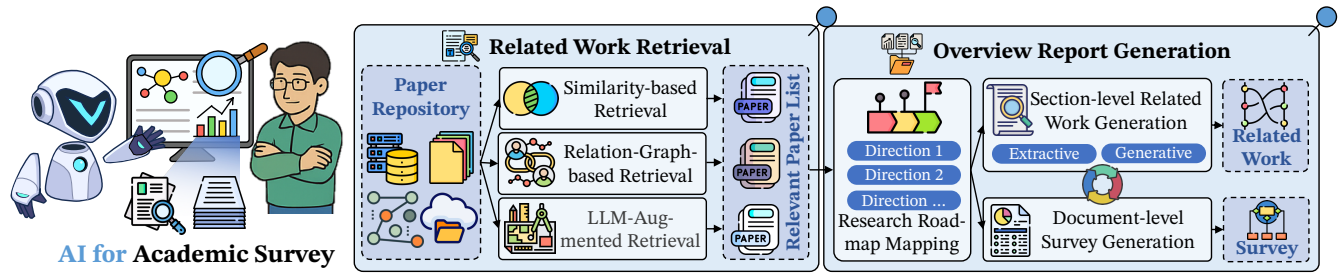


Figure 4: The two primary stages in AI-driven academic surveys: Related Work Retrieval and Overview Report Generation. Related Work Retrieval is further subdivided into Semantic-Guided Retrieval, Graph-Guided Retrieval, and LLM-Augmented Retrieval. Overview Report Generation encompasses Research Roadmap Mapping, Section-level Related Work Generation, and Document-level Survey Generation.

4.1. Related Work Retrieval

Related Work Retrieval entails AI proactively identifying foundational and novel research papers aligned with their evolving scientific objectives [52, 455, 671, 191]. Existing research divides into three paradigms:

Semantic-Guided Retrieval involves identifying relevant literature by matching the semantic representation extracted from a user query to the terms present in documents based on similarity [37, 386, 445]. In the biomedical domain, GTSLNet [512] enhances semantic-guided retrieval by utilizing a group-based keyword similarity learning network, which automatically selects clinically analogous studies. Similarly, SPLADE V2 [210] advances neural retrieval by integrating sparse lexical signals with dense expansion models, achieving significant results. Moreover, Garikaparathi et al. [224] concentrate on inspiration retrieval to provide enhanced support for idea mining.

Graph-Guided Retrieval models scholarly entities (e.g., papers, authors, citations) as a graph [74, 726, 856]. Based on the types and granularity of nodes, this search method can be categorized into three types: (1) **Author Relationship Graph** captures the connections between researchers, enabling searches based on author relationships. For example, author-relationship graphs can effectively model collaboration networks and researcher influence [360, 597, 713]. Building on this concept, Kang et al. [359] developed a “user-recommended paper” knowledge graph that traces users’ interactions with literature, enhancing recommendation transparency and trust. (2) **Paper Relationship Graph** is always constructed using citation relationships between papers to construct broader paper relationships [309]. CitationSum [516] creates a citation graph linking target papers to their references with weighted relevance scores, then uses graph contrastive learning to produce abstractive summaries. (3) **Entity Relationship Graph**: can be constructed by modeling relationships between logical entities within papers, enabling more precise retrieval. For example, Li and Ouyang [446] propose a graph-based model that automatically identifies inter-paper entity relationships, such as contrast and support, guiding the construction of a structured Related Work section. Cross-domain graph methods are increasingly used in interdisciplinary research. Further, to address complex questions such as “Which synthesis pathways enable material X to achieve optimal conductivity?”, Ye et al. [866] systematically extract entities and their relationships from the materials science literature, thus enhancing the depth of exploration.

LLM-Augmented Retrieval involves leveraging the capabilities of LLMs to improve search effectiveness and result quality by integrating them with academic retrieval systems. (1) **Single-Agent Retrieval**: The most straightforward approach is to employ a single AI model as a standalone agent to accomplish the retrieval task. For instance, Lee et al. [420] introduce the PaperWeaver framework, which places an LLM-based

agent atop a graph to enable deeper reasoning, thereby enhancing interpretability in classification and recommendation tasks. **(2) Multi-Agent Retrieval:** Beyond single-agent systems, several studies employ multiple specialized agents to simplify retrieval and increase accuracy [667, 496, 686]. LitLLMs [9] splits the automatic literature review into two subtasks: retrieval and generation. It proposes a two-phase LLM pipeline that extracts keywords from abstracts and reranks results to improve recall. Liu et al. [496] propose a multi-agent framework for full-text related-work generation, which includes a selector to choose sections to read, a reader to update shared memory, and a writer to generate the related work section. The framework uses graph-aware strategies to optimize the reading order of references. **(3) Deep Research:** Recent research has advanced this paradigm towards more autonomous “Deep Research” [576, 577], where AI agents perform the end-to-end research process, from exploration and synthesis to generating citation-rich reports [859, 178, 932]. This progress is enabled by novel agent architectures that emulate human research heuristics [824]. For instance, the PaSa agent [285] discovers literature by actively traversing citation networks. Concurrently, the retrieval strategies themselves have become more intelligent; the ExSearch framework [677] allows an agent to continuously optimize its search strategies through a self-incentivization loop, while CuriousLLM [862] employs a “curiosity-driven” mechanism where the agent actively generates questions to guide its retrieval process of knowledge graphs.

4.2. Overview Report Generation

Based on retrieved data, automated generation of structured, coherent overview reports has become essential in academic writing and AI4Research process [291]. According to the writing sequence, we need to first complete the research roadmap mapping, followed by the generation of section-level related work, and finally produce the complete document-level survey.

4.2.1. Research Roadmap Mapping

Research Roadmap Mapping refers to the process of cleaning, integrating, and depicting the developmental trajectories of a research topic by synthesizing insights from a broad corpus of literature [8, 849, 687, 802]. This methodology is crucial for enhancing the rigor and completeness of literature surveys and meta-analyses, as it enables researchers to discern emerging trends, unresolved gaps, and potential future directions more systematically [110, 69, 849]. Specifically, Zhu et al. [938] demonstrate that organizing a survey into a hierarchical structure significantly improves coherence [673].

Recently, more interactive hierarchical frameworks have also emerged. CHIME [299], for instance, refines LLM-generated structures through iterative human-AI collaboration. Similarly, Katz et al. [367] expand this to a two-tiered hierarchy, effectively organizing extensive surveys. Further, HiReview [309] illustrates the benefits of multilayered tree structures for systematic knowledge organization. Moreover, Zhuang and Kennington [945] propose a graph-based taxonomy that categorizes LLM survey papers into defined classes, outperforming fine-tuned LLMs and providing a scalable framework for organizing survey literature.

4.2.2. Section-level Related Work Generation

Section-level Related Work Generation has been regarded as a prominent research [117, 118, 580, 496, 355, 445]. Such section-level approaches are well-aligned with the actual structure of scientific papers and can effectively fulfill the requirements of the related-work-section [291].

Extractive Related Work. Early automated methods for generating the “Related Work” section involve extracting key sentences from multiple papers, which were then rewritten and combined into a coherent narrative [291]. A subsequent approach refine this by selecting papers that cited similar references to the target work and extracting relevant sentences from them [103, 788]. Further research has focused on

Methods	Model	Reference Quality		Outline Quality	Content Quality			Avg
		Input Cov.	Reference Cov.		Structure	Relevance	Coverage	
Human-Written	-	-	0.6294	87.62	-	-	-	-
AutoSurvey [802]	Claude-3-Haiku [24]	0.1153	0.2341	82.18	72.83	76.44	72.35	73.87
SurveyForge [849]	Claude-3-Haiku [24]	0.2231	0.3960	86.85	73.82	79.62	75.59	76.34
AutoSurvey [802]	GPT-4o-mini [641]	0.0665	0.2035	83.10	74.66	74.16	76.33	75.05
SurveyForge [849]	GPT-4o-mini [641]	0.2018	0.4236	86.62	77.10	76.94	77.15	77.06
SurveyForge [849]	DeepSeek-v3 [477]	0.2554	0.4553	87.42	79.20	80.17	81.07	80.15

Table 2: A comparison of document-level survey generation capabilities on SurveyBench [849] using three key Survey Assessment Metrics: Reference quality, Outline quality, and Content quality. “Input Cov.” indicates the overlap between retrieved papers and benchmark references, while “Reference Cov.” evaluates the alignment of cited references with the benchmark. Data are sourced from Yan et al. [849].

improving the organization and integration of these extracted sentences. Some methods explore optimal reference structures and sentence orderings [308, 166]. For instance, ReWoS [291] and RWS-Cit [103] build topic trees to sequence sentences, while Wang et al. [807] employ a ranking mechanism based on predicted salience probabilities to enhance the quality of the extractive related work.

Generative Related Work. Recent studies have focused on methods to structure citations and generate cohesive connecting text for entire related work sections [789, 486, 446]. These approaches generally fall into three categories: **(1) Human-Guided Generation:** This approach incorporates human input, such as keywords, short abstracts, or paper groupings, to guide the generation process and maintain focus [447, 446, 533]. For instance, Gu and Hahnloser [255] and Li et al. [448] integrate user-provided or self-extracted keywords for better related-work generation. **(2) Graph-Guided Generation:** These methods utilize citation relationships through bibliographic graphs [226, 117, 875]. Specifically, Wang et al. [807] enhance related work generation by performing random walks on heterogeneous citation graphs. Similarly, Chen et al. [118] use a graph to link references to the paper, while Li and Ouyang [446] further prompt LLMs with inter-paper features. **(3) Model-Guided Generation:** In this approach, models complete the task autonomously, without additional human input [676, 790, 603]. Guo et al. [265] and Nishimura et al. [571] treat related work generation as a summarization task with structured paragraphs and novelty statements. Additionally, Pu and Demberg [605] integrate Rhetorical Structure Theory into LoRA-based fine-tuning to identify discourse relations, and Achkar et al. [4] propose a customizable multi-stage pipeline (retrieval, citation extraction, context aggregation, polishing), further enhancing the related work generation processes.

4.2.3. Document-level Survey Generation

Document-level survey generation seeks to automate the creation of systematic literature reviews by leveraging existing research and established frameworks [814, 938, 69, 215, 825]. The detailed comparison results can be found in Table 2. For example, AutoSurvey [802] employs cue-word guidance to direct LLMs through a staged generation process. Similarly, LitLLM [7] enhances content structuring by implementing a plan-based search mechanism. SurveyX [462] strengthens logical coherence through the combination of online reference retrieval and AttributeTree preprocessing. Building on these approaches, SurveyForge [849] retrieves high-quality papers via scholar-navigating intelligences and generates survey chapters from a predefined outline, followed by iterative refinement to maintain document-level quality [696]. In contrast, STORM [673] uses multi-agent dialogue to further enhance generation performance. Beyond training-free agent management, Bio-SIEVE [642] and Susnjak et al. [731] fine-tune LLMs specifically for survey generation [403], while OpenScholar [26] offers a pipeline for training models for survey writing without relying on specialized

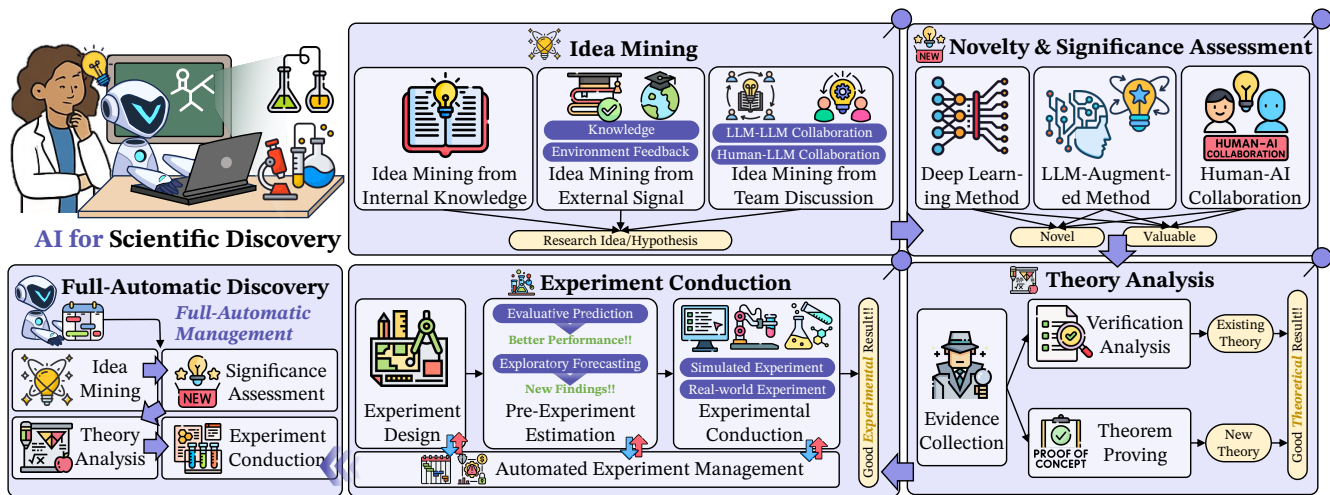


Figure 5: The AI-augmented pipeline for scientific discovery, encompassing Idea Mining, Novelty & Significance Assessment, Theory Analysis, and Experiment Conduction. Full-Automatic Discovery integrates these elements into a cohesive, end-to-end process, supporting scientific exploration and innovation.

generation architectures.

5. AI for Scientific Discovery

AI for Scientific Discovery [779, 652, 625] leverages AI to generate novel hypotheses, theories, or ideas based on existing knowledge. Its goal is to expedite the research process by automating tasks such as idea generation, novelty and significance evaluation, theoretical analysis, and experimental design. This approach not only guides new research directions but also addresses complex scientific challenges [508, 277, 284]. As shown in Figure 5, it contains five main categories: Idea Mining (§ 5.1), Novelty & Significance Assessment (§ 5.2), Theory Analysis (§ 5.3), Experiment Conduction (§ 5.4) and Full-Automatic Discovery (§ 5.5).

5.1. Idea Mining

Idea mining, also known as hypothesis generation, is crucial for producing innovative, impactful research [575, 35, 144]. Recent studies show that the LLMs exhibit strong creativity and can facilitate automated scientific discovery [500, 396, 684, 256]. A comprehensive comparison of these findings is presented in Table 3. This suggests a future where AI agents act as independent researchers. Current efforts in this domain focus on extracting ideas from various sources to foster innovation [661, 714, 474]. These methods can be broadly categorized into three main approaches:

5.1.1. Idea Mining from Internal Knowledge

Idea mining from internal knowledge leverages the latent knowledge and generative capabilities of large language models to discover novel concepts without relying on external data [237, 684, 396]. By leveraging pretrained parameters and customized prompts, researchers can extract a variety of high-quality ideas embedded within the model [540, 104]. Earlier, Meincke et al. [541] guided LLMs toward distinct “idea spaces” by adjusting decoding temperatures and applying constraint-based prompts, effectively encouraging exploration of diverse thematic trajectories. Building on these insights, Haarmann [270] conduct an empirical study where undergraduates use an interactive tool that auto-prompted GPT-4 with business-model templates,

Model	Fluency	Feasibility	Clarity	Originality	Flexibility	Average
DeepSeek-R1 [263]	6.63	6.52	8.10	7.84	6.83	7.18
Deepseek-R1-Distill-Qwen-32B [263]	7.06	6.08	7.43	7.13	6.62	6.86
Deepseek-R1-Distill-Llama-70B [263]	6.66	6.07	7.43	6.98	6.41	6.71
Claude-3.7-Sonnet [24]	7.80	5.46	7.61	7.81	6.92	7.12
Claude-3.5-Sonnet [24]	6.90	5.42	7.85	7.83	6.62	6.92
Claude-3.5-Haiku [24]	5.61	5.05	7.40	7.72	6.08	6.37
Claude-3-Opus [24]	5.74	5.66	7.72	6.66	6.04	6.36
Gemini-2.0-Flash-Exp [744]	7.30	6.02	7.84	7.37	6.83	7.18
Gemini-2.0-Flash-Thinking-Exp [744]	7.38	6.05	7.69	7.35	6.83	7.06
Gemini-2.0-Pro-Exp [744]	6.84	5.88	7.90	7.76	6.75	7.03
Gemini-Pro-1.5 [745]	6.68	5.92	7.75	7.33	6.58	6.85
GPT-4o [3]	6.12	5.58	7.74	7.64	6.38	6.69
o1-mini [334]	5.89	6.20	7.77	7.09	6.33	6.66
o1 [334]	6.23	5.88	7.42	7.23	6.29	6.61
o3-mini [334]	5.57	5.91	7.43	7.45	6.21	6.51
o3-mini-high [334]	5.76	5.82	7.62	6.95	6.17	6.47
GPT-4o-mini [3]	5.28	5.86	7.45	6.67	6.00	6.25
Llama-3.1-405B-Instruct [250]	6.57	5.56	7.48	7.18	6.33	6.62
Llama-3.1-70b-Instruct [250]	6.71	5.49	7.34	7.16	6.38	6.62
QwQ-32B [751]	6.45	6.35	7.98	7.77	6.75	7.06
QwQ-32B-Preview [750]	7.49	6.10	7.46	6.87	6.71	6.93
Qwen-2.5-72B-Instruct [851]	6.17	5.99	7.72	6.91	6.29	6.62
Qwen-2.5-7B-Instruct [851]	6.66	6.02	7.17	6.34	6.17	6.47
Mistral-Small [342]	7.36	5.97	7.36	6.98	6.62	6.86
Mistral-Large [342]	6.68	6.06	7.69	7.01	6.50	6.79
Nova-Pro-v1 [331]	6.45	6.19	7.41	6.59	6.25	6.58
Nova-Lite-v1 [331]	4.51	6.06	7.38	6.60	5.71	6.05
Phi-4 [1]	6.58	5.80	7.57	7.24	6.42	6.72
Gemma-2-27b-IT [746]	7.18	5.50	7.36	6.86	6.38	6.65
Grok-2 [836]	5.76	5.82	7.62	6.95	6.17	6.47

Table 3: Results from the Liveideabench benchmark [655] across five key dimensions: Fluency, Feasibility, Clarity, Originality, and Flexibility. Data is sourced from Ruan et al. [655]. The **bolded** contents indicate the highest performance for each metric.

resulting in ideas with higher novelty and feasibility, all without extensive innovation training. Additionally, Liu et al. [495] demonstrate that injecting metadata into the LLM-based ideation process and applying automated validation during selection significantly increased idea feasibility and overall quality in climate-negotiation experiments. Furthermore, Chen et al. [111] model the process of inference-time learning and reasoning as a circuit, enhancing the idea-mining ability of the model through various voltage-enhancing techniques.

5.1.2. Idea Mining from External Signal

LLMs in research workflows can leverage not only internal parameterized knowledge, but also external signals to generate more novel, feasible, and contextually relevant hypotheses and ideas. By incorporating structured knowledge repositories or experimental feedback, these approaches extend beyond purely internal reasoning. We categorize them into two types:

Idea Mining from External Knowledge involves supplying AI with curated academic data, such as publication metadata, citation networks, or knowledge graphs, to drive idea mining. Integrating up-to-date, domain-specific information ensures that generated hypotheses align with the latest developments in the field [483, 440, 575]. Early efforts, limited by the capacity of earlier language models, focus on predicting relationships between concepts to generate classical “A+B” ideas [288, 385]. With recent advancements in language modeling [616, 919], attention has shifted to utilizing LLMs to explore ideas from scholarly data [526, 483]. To support this, researchers have proposed various strategies for organizing literature to optimize knowledge extraction and mining. These mainly include ternary representations [793], chained structures [433], comprehensive databases [796], and knowledge graphs [79, 257, 258, 221]. Furthermore, efforts have focused on refining the knowledge injection process in idea mining. For example, Gu et al. [256] propose a framework with (1) a generalized retrieval system for cross-domain knowledge discovery and (2) a structured combinatorial process for improved idea mining.

Idea Mining from External Environment Feedback involves treating idea mining as an interactive loop that incorporates feedback from experimental or simulated environments [36, 606]. Static document mining, by contrast, often overlooks the complexities of the real world, limiting its potential for innovation. These methods enable AI systems to propose experiments, receive outcome data, and refine subsequent ideas, thereby mimicking the research cycle of design, execution, and analysis [29, 12]. In this domain, researchers primarily utilize multi-agent-based autonomous research systems, integrating idea mining agents with experiment conduction agents [666, 665, 330]. Furthermore, researchers have successfully extended idea mining to various experimental disciplines, including chemistry [518], materials science [566], biology [151], medicine [732, 249], and machine learning [320].

5.1.3. Idea Mining from Team discussion

Idea mining from team discussion encompasses approaches that simulate or facilitate collaborative brainstorming among multiple agents, either purely algorithmic or involving human participants, leveraging iterative critique, background knowledge retrieval, and facet recombination to generate richer, more diverse idea portfolios than single-agent pipelines.

AI-AI Collaboration improves scientific ideation by refining hypotheses, critiquing proposals, and integrating external knowledge (also referred to as multi-agent collaboration) [732, 502]. We categorize current approaches into two mainstreams: (1) **Feedback-guided Mining** involves agents exchanging critiques at various research stages to refine hypotheses through iterative feedback. Some studies introduce feedback loops across idea mining, experimental design, and result interpretation to optimize performance [935, 860, 695], while others refine hypotheses using earlier outputs [306]. These methods integrate peer review [507], direct critiques of hypotheses [36], and evaluations of experimental results [519, 877]. (2) **Team-Discussion-guided Mining** assembles multiple agents with distinct roles to simulate human research team dynamics [607, 568, 231]. Specifically, Su et al. [717] create a virtual research team (VirSci) where agents iteratively propose and critique ideas, producing more novel concepts than single-agent prompts by leveraging an expanding idea archive [892, 507]. Yang et al. [861] has developed a multi-intelligentsia framework, MOOSE-Chem, based on LLMs, specialized in scientific hypothesis discovery in chemistry, which can perform the functions of retrieving inspiration and generating hypotheses based on research contexts. Moreover, Li et al. [433] introduce the Chain-of-Ideas (CoI) agent, which organizes literature into a sequential chain, mirroring a topic’s historical progression. This method generates outputs of similar quality to small research teams with minimal costs. Lagzian et al. [402] further enhance diversity and novelty via inference-time multi-view brainstorming.

Human-AI Collaboration means the process where a human researcher guides an LLM’s exploration by selecting and curating intermediate artifacts, which the model then recombines and refines [567, 566]. For instance, Radensky et al. [623] introduce Scideator, a system that enables researchers to select various facets, such as the problem statement, methodology, and dataset, from existing papers. The LLM subsequently recombines these facets to generate novel candidate ideas, significantly improving the idea qualities. Similarly, Garikaparthi et al. [225] present IRIS, an interactive research ideation system that facilitates human-AI collaboration by validating research motivations and synthesizing methodological suggestions in response to researcher queries. However, the research findings [394] show that, although LLM assistance can yield short-term boosts in creativity during supported tasks, it may inadvertently hamper users’ independent creative performance when working unassisted, thereby raising concerns about its long-term effects on human creativity and cognitive abilities.

5.2. Novelty & Significance Assessment

Novelty & Significance Assessment focuses on AI methods that evaluate the originality and impact of ideas and scholarly papers [351, 684, 370]. The field predominantly employs three approaches: **(1) Traditional Methods:** Initially, models are trained to classify or regressively assess novelty and significance [177, 329, 801]. For instance, Singh and Chakrabarti [689] propose SAPPhIRE that utilizes the causality ontology to quantify novelty in design problems, measuring textual similarity at multiple abstraction levels against historical works. Additionally, Wang [778] introduce “surprise” as an alternative measure of novelty, comparing a paper’s word distribution to a language model’s representation of scholarly discourse. This approach aligns with scientific intuition (face validity) and shows a correlation with expert judgments (construct validity). **(2) LLM-Augmented Methods:** With the significant development of LLMs, a series of works try to integrate LLMs for better novelty and significance assessment [831]. Typically, Feng et al. [204] propose GraphEval, a lightweight, graph-based LLM framework for reasoning evaluation, that prompts a small-scale LLM to decompose complex reasoning processes into easily interpretable “viewpoint” nodes, thereby enhancing the robustness of reasoning assessment. **(3) Human-AI Collaboration Methods:** Unfortunately, purely LLM-augmented assessments of novelty may overestimate creativity [93, 412] and lead to homogenization effects without human input [21, 936]. As a result, there is growing interest in human-AI collaboration for novelty assessment, with several works [27, 499, 583] advocating for the integration of human-guided ideation alongside LLM-based workflows.

5.3. Theory Analysis

Any scientific idea or hypothesis must be rigorously evaluated to confirm its validity. Theory analysis involves using AI methods to determine whether a hypothesis aligns with established scientific principles. AI applications in theory analysis can be divided into three main components:

5.3.1. Scientific Claim Formalization

Scientific claim formalization converts natural-language assertions into structured representations for systematic verification [633, 397, 547]. Early approaches relied on template-based methods [247]. For example, Heger et al. [286] describe a pipeline that converts complex hypotheses into machine-readable templates. Subsequent works focus on incorporating LLMs to refine these templates. Ganguly et al. [218] propose a PCFG-based framework to address common LLM failure modes, while Valsci [182] automates the conversion of natural-language claims into templated queries for LLM-driven verification. More recently, Yan et al. [850] suggest that integrating text, images, and other modalities in multimodal LLMs provides richer structured representations, facilitating cross-domain reasoning.

5.3.2. Scientific Evidence Collection

Scientific evidence collection involves systematically identifying, retrieving, and curating data sources to support or challenge research claims [586, 768, 368]. Previous studies have focused on methods for evaluating and improving the quality of retrieved sources [765] and optimizing retrieval configurations [766]. Additionally, strategies have emerged to address incomplete or faulty evidence, including techniques for detecting missing information [239] and understanding the causes of retrieval errors [835, 240]. More recent efforts have integrated LLMs with retrieval systems to enhance the accuracy of evidence retrieval and verification. For instance, SciClaims [578] combines claim extraction, evidence retrieval, and verification into a single LLM-powered pipeline, streamlining the entire process. Similarly, Alvarez et al. [17] and Wang et al. [795] extend retrieval-augmented generation by producing structured query representations and retrieving corroborating or refuting evidence in a single step.

5.3.3. Scientific Verification Analysis

Scientific verification analysis plays a critical role in AI-driven theoretical qualitative studies by assessing the logical coherence [390, 460], factual consistency [554, 363, 91, 75], and robustness [335] of claims based on existing evidence. Research underscores the importance of domain expertise for accurate and reliable verification [160, 16, 50, 834]. To mitigate errors and enhance interpretability, some frameworks adopt human-like, stepwise pipelines [376, 155, 826, 30]. For instance, HiSS [903] and ProToCo [884] employ multiple cueing to validate each substatement, improving reliability. Other methods integrate verification with experimental results to boost transparency and interpretability [387, 585, 187, 900]. GX-Chen et al. [269] show that LLMs inherit reasoning heuristics from training data, leading to cognitive biases. To mitigate this, they propose an inference-time-scaling sampling procedure that reduces implicit causal assumptions and aligns the model’s reasoning with causal rigor. More recently, Ku et al. [390] introduced the task of generating coherent visual explanations and demonstrated that combining agents with Manim animations to produce long-form theorem explanation videos (over five minutes) results in more effective visual explanations.

5.3.4. Theorem Proving

Theorem proving involves the development of algorithms and models, often incorporating generative language models, to autonomously generate and verify formal mathematical proofs [454, 212, 937, 853]. Early methods [776, 407] introduce dynamic tree proof search techniques and integrate retrieval algorithms with language models for theorem proving [599]. However, retrieval algorithms tend to prioritize trivial intermediate conjectures, resulting in poor performance [777]. To overcome this, some researchers have introduced novel approaches that replace retrieval algorithms entirely [341, 340]. LEGO-Prover [775] employs Growing Libraries to enhance LLM reasoning, while Zhao et al. [922] suggest Subgoal-based Demonstration Learning for more effective theorem proving. Additionally, Lean Copilot [703] and Lean-STaR [468] leverage the Lean programming language and theorem prover to enable improved human-AI collaboration in proof completion. Recent studies have focused on fine-tuning specialized proving LLMs [208]. For example, MUSTARD [323] and DeepSeek-Prover [841] aim to generate high-quality synthetic data to fine-tune models and improve theorem proving.

5.4. Scientific Experiment Conduction

Automatic Scientific Experiment Conduction leverages AI to design, conduct, and analyze scientific studies autonomously, aiming to automate the entire process, from hypothesis formulation to data interpretation. This automation seeks to accelerate research and improve reproducibility [112, 383, 49, 764]. However, Zhu et al. [941] highlight a critical challenge: AI scientists currently lack the validation capabilities needed

for rigorous experimentation and high-quality manuscript production. Without these essential competencies, such platforms cannot succeed.

5.4.1. Experiment Design

Experimental design is vital for efficiency and provides the foundation for AI-assisted experiment conduction methods [764, 151, 249, 231]. Evidence shows that systematic design plays a central role in automating and enhancing experimental processes [566, 518].

Semi-Automatic Experiment Design involves the creation of experimental plans through human-AI collaboration [566, 518]. Arlt et al. [25] present a transformer-based framework that autonomously generates quantum experiment protocols and uncovers state preparation principles. Huang et al. [322] integrate deep learning with multi-objective optimization to design polymer sequences with both high thermal conductivity and synthetic feasibility, validating their results through molecular dynamics. Liang [459] apply a variational autoencoder combined with reinforcement learning to enhance the design and efficiency of parameters for cultural creative products. Craig [148] propose a human-AI collaboration framework based on experimental design, case-based reasoning, and a note-taking system, offering scientists a structured LLM tool with transparent documentation, resulting in verifiable experimental designs and knowledge integration.

Full-Automatic Experiment Design refers to the application of agent-centric methods for the automatic scheduling of scientific experiments [330, 12, 233]. Platforms such as The AI Scientist [507] and Agent Laboratory [666] continuously refine experimental protocols by incorporating new data in real-time [732, 36, 665]. In a significant development, Liu et al. [485] proposed an end-to-end generative-agent framework that enables fully autonomous planning, spanning from literature review to protocol iteration, without the need for human intervention. Additionally, Roohani et al. [648] introduced a biodiscovery agent capable of designing, evaluating, and optimizing gene-perturbation experiments. This system has shown superior performance over traditional Bayesian methods, especially in targeting non-essential genes.

5.4.2. Pre-Experiment Estimation

Pre-experiment prediction leverages AI to forecast experimental outcomes, aiming to improve research efficiency and accuracy. This process can be divided into two categories:

Evaluative Prediction predicts quantitative values or trends of experimental outcomes, such as estimating drug concentration effects, determining whether a compound affects cellular activity, and assessing protocol feasibility [151]. **(1) Deep-Learning Methods:** With the rise of deep learning, hierarchical prediction models have emerged [833]. Li et al. [437] incorporated physical equations to predict pharmacokinetic parameters, reducing data requirements and enhancing noise robustness. More recently, Li et al. [438] proposed a dual-matching framework, combining hierarchical molecular alignment with meta-learning, which showed significant improvements in drug feature estimation. **(2) LLM-Augmented Methods:** More recently, with the advent of LLM capabilities, Zhang et al. [901] demonstrate successful LLM-assisted evaluative prediction, a method that has been further extended in subsequent studies. Notably, Luo et al. [515] integrated BrainGPT into neuroscience literature retrieval, outperforming domain experts in evaluating experimental estimation. Wen et al. [817] developed a system combining fine-tuned GPT-4.1 with a paper retrieval agent, which outperformed 25 human experts in evaluating experimental predictions.

Exploratory Forecasting utilizes AI to predict experimental outcomes, generate new compounds, design reaction pathways, and propose combinatorial schemes to drive scientific discovery [518, 494]. Several studies have applied deep generative models for chemical-space forecasting [244, 163]. Seo et al. [669] introduce a framework that uses graph diffusion modeling to predict ingredient-chemical molecule interactions,

enabling innovative pairing exploration. Furthermore, based on a massive computation model, DeepMind’s GNoME [41] predicts approximately 380,000 stable material structures, demonstrating AI’s potential in materials discovery. Recently, multi-turn interactive methods have also been developed to improve exploratory forecasting [249]. For instance, Zhang et al. [901] and Swanson et al. [732] present platforms that integrate multi-agent debates to better forecast experimental performance and foster idea exchange, advancing the discovery of new method variants.

5.4.3. Experiment Management

The integration of machine learning and robotics in AI-driven experiment management enables hypothesis generation, high-throughput experimentation, and iterative procedure refinement without human intervention [44, 672, 923, 23]. These paradigms, also named as “self-driving laboratories” [76, 87, 281], promise accelerated discoveries [198] in biology, chemistry, and materials science [400, 89, 197].

Open-Loop Management involves experimental management without human oversight [772]. Hysmith et al. [326] explore human-AI collaboration, emphasizing the interoperability of robots, predictive models, and data pipelines. In bioprocessing, Zournas et al. [950] combine active learning with a semi-automated Design-Build-Test-Learn cycle to optimize microbial media, showing that higher NaCl levels significantly improve metabolite yield and process efficiency. Google DeepMind and BioNTech [207] have introduced an AI-driven laboratory assistant to autonomously design protocols and predict outcomes, aiming to enhance research efficiency in the medical, energy, and educational sectors. Reports indicate that such systems could reduce the traditional 20-year, \$100 million timeline for materials discovery to just months. The U.S. government is supporting these efforts through strategic funding initiatives [33].

Close-Loop Management entails fully autonomous experimental management without human intervention [772]. The Functional Genomics Explorer [379] is a landmark in this area, being the first fully autonomous research platform that generates hypotheses, designs experiments, and validates results. MacLeod et al. [524] describe a robotic system that formulates, deposits, and characterizes thin films using model-based optimization to enhance charge transport. AI-driven optimization algorithms are transforming experimental workflows, as seen in closed-loop Bayesian optimization methods for chemical and materials discovery [756]. Knox et al. [381] apply multi-objective optimization to polymer nanoparticle synthesis, optimizing size, dispersity, and functionality.

5.4.4. Experimental Conduction

Experimental conduction refers to the application of AI techniques in executing and managing scientific experiments. This process is essential for automating workflows, ensuring experiments are carried out efficiently and accurately. The primary aim of experimental conduction is to reduce human involvement in the experimental process. It can be further divided into two categories:

Automated Machine Learning Experiment Conduction uses AI to streamline the design, training, and evaluation of ML models, reducing dependence on human expertise by covering the entire pipeline from preprocessing to hyperparameter optimization [897, 758, 253, 911, 574, 927]. Typically, Wang et al. [799] present a community-driven sandbox allowing agents to write code, browse the web, and coordinate through an event-stream API. For Kaggle challenges, Li et al. [457] propose an iterative, collaborative multi-agent system that incorporates debugging and unit testing across the competition pipeline. AIDE [174] utilizes a tree-search loop to generate, evaluate, and refine solutions, achieving a bronze medal in Kaggle competitions. At the research level, Li et al. [441] formalize a three-phase LLM agent workflow, idea generation, implementation, and execution, to automate experiments. Zhao et al. [921] and Liu et al. [481]

Models	Without Knowledge				With Knowledge			
	SR	CBS	VER	Cost ↓	SR	CBS	VER	Cost ↓
<i>Direct Prompting</i>								
Llama-3.1-Instruct-70B [250]	5.9	81.5	29.4	0.001	4.9	82.1	27.5	0.001
Llama-3.1-Instruct-405B [250]	3.9	79.4	35.3	0.010	2.9	81.5	25.5	0.001
Mistral-Large-2 [342]	13.7	82.3	47.1	0.009	16.7	84.7	39.2	0.001
GPT-4o [3]	11.8	82.6	52.9	0.011	10.8	83.8	41.2	0.016
Claude-3.5-Sonnet [24]	17.7	83.6	51.0	0.017	21.6	85.4	41.2	0.016
o1-preview	34.3	87.1	70.6	0.221	31.4	87.4	63.7	0.236
<i>OpenHands CodeAct [800]</i>								
Llama-3.1-Instruct-70B [250]	6.9	63.5	30.4	0.145	2.9	65.7	25.5	0.252
Llama-3.1-Instruct-405B [250]	5.9	65.3	32.0	0.383	8.3	71.4	58.0	0.384
Mistral-Large-2 [342]	9.8	72.5	53.9	0.513	13.7	78.8	50.0	0.759
GPT-4o [3]	19.6	83.4	87.5	0.803	27.5	86.3	73.5	1.094
Claude-3.5-Sonnet [24]	21.6	83.6	87.3	0.958	24.5	85.1	88.2	0.900
o1-preview [334]	33.4	86.2	87.0	0.999	35.3	88.4	91.5	0.913
<i>Self-Debug [116]</i>								
Llama-3.1-Instruct-70B	13.7	82.7	40.4	0.007	16.7	83.5	73.5	0.005
Llama-3.1-Instruct-405B	16.7	80.0	35.3	0.006	23.6	79.4	40.4	0.004
Mistral-Large-2 [342]	23.5	85.1	78.4	0.007	26.5	86.7	84.3	0.006
GPT-4o [3]	22.6	84.4	84.3	0.024	33.4	87.1	86.3	0.037
Claude-3.5-Sonnet [24]	32.4	86.4	92.2	0.026	34.5	87.1	86.3	0.015
o1-preview [334]	42.2	88.4	92.0	0.636	41.2	88.9	91.2	0.713

Table 4: Full-automatic discovery capability comparison on ScienceAgentBench [126]. The data presented are derived from Chen et al. [126]. The **bolded** contents indicate the highest performance for each metric.

present a multi-agent method for extracting model variables from scientific texts, significantly improving experimental reproduction accuracy.

Real-world Experimental Simulation & Conduction. Recent advancements in the planning and reasoning capabilities of LLMs have led to their use in simulating experimental results [357, 882, 414] and even direct conduct real-world experiments [122, 656, 200]. Real-world experimental simulation & conduction generally employ four strategies: *(1) Self-Improvement:* Models iteratively refine their performance based on feedback [311, 680, 908, 593, 877]. For example, Siddiqui et al. [685] enhance functional approximation through iterative knowledge application. Further refinement occurs through analytical insights [441, 734, 36] and hyperparameter tuning [566, 492, 893]. *(2) Multi-Agent Interaction:* Models simulate collaborative research teams by assigning roles such as experimenter, analyst, or critic [36, 231, 710, 305, 699]. For instance, MechAgent [565], Researchcodeagent [216] and The AI Scientist [507, 848] automate experiments through multi-agent collaboration, with LLMs acting as proxies in fields like computer science [507, 848], social science [550, 530], and physical science [898]. *(3) External Tool Integration:* Researchers enhance model capabilities by linking them to databases, APIs, and other tools during experiments [276, 618, 663]. For example, Boiko et al. [68] integrate internet search, code execution, and automation into a GPT-4 system. Studies like ChemCrow [518] and Crispr-GPT [314] support chemistry and gene editing experiments through massive specialized toolchains. *(4) Specific Fine-Tuning:* A growing body of work explores the fine-tuning of specific models to improve experimental simulations. For instance, Cui et al. [150] present a transformer-based model trained on large single-cell transcriptomic datasets, achieving state-of-the-art accuracy in cell-type annotation and in silico perturbation response [489].

5.4.5. Experimental Analysis

Experimental Analysis involves systematically testing hypotheses, evaluating models, or validating theoretical assumptions to draw meaningful conclusions. This process encompasses three main sub-processes:

Automated Evaluation Metrics refer to systems like AutoML that automatically generate model learning curves, parameter sensitivity analysis graphs, and other evaluation tools to assess model performance [5]. For instance, AutoML platforms assist researchers by automatically producing learning curves and sensitivity analysis graphs to better understand model behavior [42, 40].

Theoretical Consistency Analysis ensures that the theoretical methods align with the experimental implementations [481]. AutoReproduce [921] uses a large language model to create a multi-intelligent body system, enabling automatic comprehension, code reproduction, and execution verification of experiments in scientific papers. This process completes the consistency analysis between theoretical methods and experimental outcomes.

Exploratory Analysis is essential for investigating and understanding datasets through statistical and visualization techniques to identify patterns, spot anomalies, test assumptions, or validate hypotheses [100, 719]. This process extends the capabilities of language models for data exploration in structured formats. For example, Xing et al. [842] utilizes a generator-validator fine-tuning approach to enable language models to specialize in parsing tabular data, improving table structure inference and summarization. Additionally, Bian et al. [62] developed HeLM, which facilitates high-quality natural language summarization of table content, aiding in the generation of conclusions.

5.5. Full-Automatic Discovery

Full-automatic discovery refers to the ability to close the loop of the scientific process, from hypothesis generation and experimental design to autonomous execution, result analysis, and iterative feedback, powered by end-to-end artificial intelligence [529]. A comprehensive comparison of the results is presented in Table 4. Advances in laboratory automation and closed-loop assistants are driving fully automated discovery toward greater reliability, innovation, and faster iteration through multi-agent systems [558, 561, 749, 330]. For example, Lu et al. [507] and Yuan et al. [877] use literature mining to rank research topics, employ an “anomaly-guided” code-synthesis framework to generate and debug experimental scripts, and feed results back into the ideation module to iteratively refine hypotheses [666, 665, 538]. Kon et al. [382] introduce rigor through three modules: intra-agent rigor for reliability, inter-agent rigor for systematic control, and an experimental knowledge module for interpretability, addressing issues of insufficient rigor and overstated claims. Li et al. [451] extend this approach to data-driven discovery, enhancing exploration diversity. Further, Zochi [12] is developed as an AI-driven system for end-to-end scientific discovery, demonstrating its comprehensive capabilities across the research lifecycle. Papers generated through Zochi have even been accepted by ACL 2025.

6. AI for Academic Writing

AI for Academic writing involves the use of AI techniques to assist researchers or generate from scratch in drafting, editing, and formatting scientific manuscripts [371]. With the development of deeper interaction between human and LLMs, human and LLMs are quickly shaping each other’s better writing habits [227, 86, 933]. As shown in Figure 6, it contains two main categories: Semi-Automatic Academic Writing (§ 6.1) and Full-Automatic Academic Writing (§ 6.2).

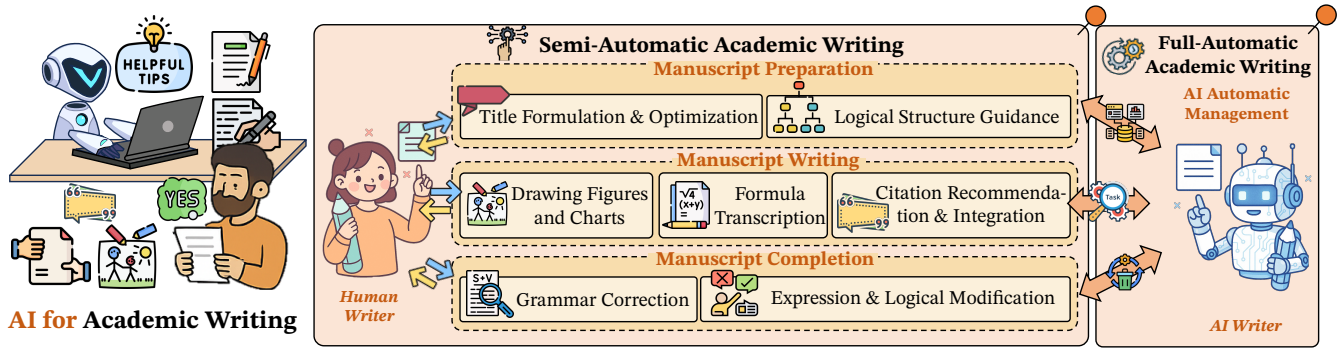


Figure 6: The main paradigms of AI for Academic Writing. It can be divided into two main categories: Semi-Automatic Academic Writing and Full-Automatic Academic Writing. Specifically, Semi-Automatic Academic Writing encompasses Manuscript Preparation, Manuscript Writing, and Manuscript Completion.

6.1. Semi-Automatic Academic Writing

Semi-automatic academic writing involves the use of AI tools to assist researchers in drafting and editing scientific manuscripts, requiring human input and oversight. This approach aims to enhance the quality and efficiency of scientific writing by offering AI-generated suggestions, corrections, and formatting assistance. Semi-automatic academic writing can be categorized into three phases:

6.1.1. Assistance During Manuscript Preparation

Assistance During Manuscript Preparation refers to the support provided by models and tools throughout the manuscript creation process. This includes generating and refining titles, guiding the overall structure, and ensuring content coherence, all aimed at improving the clarity, quality, and readiness for submission.

Title Formulation and Optimization involves using models to generate multiple title candidates and selecting the most suitable one [72, 63]. For example, given a research topic like “new energy battery materials”, a model generates 5-10 titles with varying focuses, which are evaluated based on novelty, complexity, and potential impact, helping the author select the best option [63]. To further improve the title quality, Rehman et al. [638] fine-tune PEGASUS-large and use GPT-3.5-turbo (zero-shot) to generate titles from abstracts. Au et al. [31] enhance title quality by incorporating user preferences for more coherent and personalized titles.

Overall Logical Structure Guidance involves providing the model with section and subsection headings, as well as paragraph outlines, to evaluate the logical flow, ensure the avoidance of repetition, correct ordering errors, and identify any missing elements [279]. Sun et al. [722] propose a multi-stage workflow for assessing paper structure. Their rubric emphasizes the importance of section completeness and content cohesion, ensuring that headings and paragraphs adhere to established formatting standards.

6.1.2. Assistance During Manuscript Writing

Semi-automatic AI writing involves a collaborative process where humans create the primary content, while AI contributes supplementary elements. In this collaboration, AI assists with tasks that are secondary to the main content [466, 127, 846, 701, 564]. This process can be divided into three primary tasks:

Drawing Figures and Charts serves as an effective means of conveying experimental data and analysis [55]. Recent advancements in AI research on automatic scientific figure generation have led to significant

breakthroughs [894, 293, 747, 635]. Rodriguez et al. [643] introduce the FigGen model, which maps textual descriptions to complex academic figures with high fidelity. Beyond direct image generation, several studies have explored programming scientific diagrams using Python [920, 913, 131, 551], SVG [645], or tikz [55, 54] for better figure quality. However, figures alone cannot fully convey results; captions are essential for ensuring that readers understand each figure [90, 507, 867]. To address this, Hsu et al. [300] develop SciCapenter, an interactive system that generates multiple caption candidates, scores them, and quality-checks each, assisting authors in selecting the most optimal phrasing. MLBCAP [375] utilize multiple LLMs to collaborate on chart-based title generation. Additionally, frameworks like AI Scientist [507, 848] can identify key data from experimental logs, generate figures with captions, and integrate them into authoring tools, advancing end-to-end AI-driven scientific visualization.

Formula Transcription need to digitize extensive mathematical formulas and tables in academic and instructional materials, driven the development of automated tools that convert handwritten or image-based expressions into editable LaTeX. Specifically, Sundararaj et al. [729] employ a Vision Transformer (ViT) to transcribe these expressions into LaTeX, improving accuracy and reducing manual proofreading. Vrečar et al. [767] introduce a semi-automated tool for semantic annotation, enhancing the accessibility and interoperability of mathematical symbols in LaTeX. Jiang et al. [343] propose the iterative refinement framework, which generates an initial LaTeX draft, compares it to the source image for feedback, and iteratively corrects errors, thereby minimizing manual verification and facilitating better transcription.

Citation Recommendation & Integration has emerged as a key research area in academic writing, focusing on retrieving and incorporating relevant literature into documents to enhance writing efficiency and citation accuracy [906, 525, 813, 449, 14, 465]. Earlier, Ma et al. [521] introduce a temporal preference model for ranking citations, setting the stage for subsequent research considering time factors [653]. In generative models, Çelik and Tekir [92] develop CiteBART, which masks citation markers in context and reconstructs them, enabling zero-shot citation generation. More recently, Wang et al. [808] introduce ScholarCopilot, which generates special retrieval tokens and dynamically queries literature databases to embed references in real time. He et al. [285] propose PaSa, an advanced paper-search agent powered by large language models. PaSa autonomously invokes search tools, reviews manuscripts, and selects relevant references, achieving results that surpass even Google + GPT-4o in complex scholarly queries.

6.1.3. Assistance After Manuscript Completion

After completing the paper, the author typically needs to refine its quality, focusing on language accuracy and logical coherence. At this stage, support tools can assist in grammar correction, expression and logical revision, ensuring the paper is clear, fluent, and logically sound.

Grammar Correction means the model proofreads each paragraph, identifies spelling errors, improper punctuation, repetitive phrasing, and character-encoding issues, and provides corresponding revision suggestions [228, 854, 804, 721]. Specifically, Wang et al. [805] propose a synthetic data construction method based on contextual augmentation, which can ensure an efficient augmentation of the original data with a more consistent error distribution. Wang et al. [783] propose an integrate automated writing evaluation system with grammatical error correction to support L2 essay writers by providing immediate feedback, offering targeted guidance to improve grammar and coherence, reducing manual grading efforts. Further, Zheng and Zhang [930] present a Transformer-based feedback framework that generates real-time suggestions on grammar, vocabulary, sentence structure, and logical coherence for non-native English writers. Its modular design and dynamic parameter adjustment enable personalized learning paths while ensuring low-latency feedback and differential privacy.

Expression & Logical Revision highlights AI systems' role in refining scientific manuscripts post-initial draft, focusing on expression and logic [930]. (1) **Self-guided Revision** involves AI autonomously analyzing drafts and suggesting edits to improve language, cohesion, and structure [702, 196]. Ito et al. [333] propose sentence-level edits, adjusting or rewriting sentences based on draft content. Additionally, Botha et al. [73] use revision histories to segment and rewrite the text, further enhancing revision quality. (2) **Human-guided Revision** refers to interactive systems where users provide specific instructions or highlight sections for the AI to modify, forming a collaborative editing process [582, 229, 416, 556]. Faltings et al. [188] develop an interactive editor that responds to user commands. Wordcraft [143] supports few-shot learning and dialogue for interaction. However, these methods struggle to capture the diversity and iterative nature of revision. XtraGPT [106] provides open-source LLMs for context-aware, instruction-guided revisions, addressing surface-level and section-level coherence. (3) **Human-in-loop Revision** emphasizes a cyclical workflow combining AI suggestions, human evaluation, and document updates through multiple optimization loops [202, 328, 742, 557]. Du et al. [179] propose a human-in-the-loop system that integrates model-generated edits, user feedback, and document updates for high-quality revisions [761]. Lin [476] shows that human-AI frameworks improve collaborative efficiency. Wen et al. [816] develop OverleafCopilot, a browser extension integrating LLMs into Overleaf for real-time suggestions, automatic rewriting, translation, and prompt sharing through PromptGenius to enhance LaTeX writing.

6.2. Full-Automatic Academic Writing

Full-automatic academic writing refers to using AI to generate complete scientific manuscripts without human intervention. This process spans drafting, formatting, and producing high-quality papers ready for submission, effectively removing the need for human input in manuscript preparation [786, 844]. Recent research has primarily adopted multi-agent, modular designs with self-feedback mechanisms for iterative refinement. The AI Scientist [507] treats writing and reviewing as pipeline modules: by simulating peer review and providing score-based feedback, it refines drafts. Yamada et al. [848] extend this system by incorporating vision-language model feedback loops to improve both content and figure presentation. Agent Laboratory [666, 665] employs a paper-solver module with role-based agents that simulate lab workflows, evaluate drafts against NeurIPS criteria, and iteratively enhance them. Zochi [12] uses a multi-agent architecture for initial draft generation, combining automated review with self-feedback for further polishing. Despite these successes, including some papers passing human peer review, no system has yet fully eliminated human editing, especially regarding correct citation use [330, 12, 666].

7. AI for Academic Peer Reviewing

Peer reviewing plays a crucial role in enhancing the quality of academic papers. However, it is often hindered by delays, time demands, and growing academic workloads [469, 384, 753, 946]. To address these challenges and improve dissertation quality, researchers are exploring the integration of AI into the review process [880, 490, 573, 399, 523]. As shown in Figure 7, it contains three main categories: Pre-Review (§ 7.1), In-Review (§ 7.2) and Post-Review (§ 7.3).

7.1. Pre-Review

In this phase, editors or track chairs are tasked with preliminary scoring, identifying the manuscript's subject domain, and assigning appropriate reviewers to ensure review quality and prevent conflicts of interest.

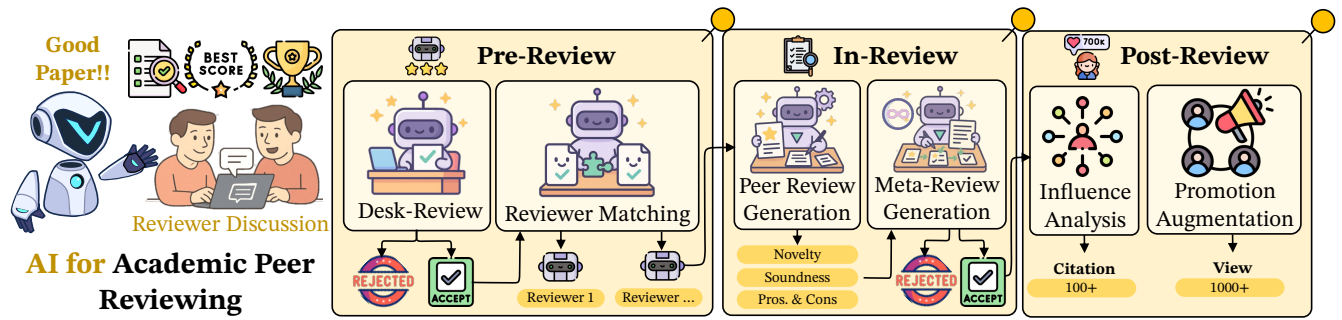


Figure 7: The primary process pipelines in AI for Academic Peer Review, encompassing three key stages: (1) Pre-Review, including Desk-Review and Reviewer Matching to ensure higher quality and more efficient evaluations; (2) In-Review, comprising Peer Review and Meta-Review, aimed at providing comprehensive scholar feedback and evaluation; and (3) Post-Review, featuring Influence Analysis and Promotion Augmentation, designed to assess the impact of the review process and improve the dissemination of scholarly work.

7.1.1. Desk-Review

As manuscript submissions to academic journals increase, editorial offices face heavier workloads during the desk-review stage. To address this, major publishers have introduced AI-driven tools, such as automated keyword extraction, topic matching, and preliminary scoring, to enhance efficiency, shorten turnaround times, and reduce manual screening [176, 426, 195]. For example, Elsevier’s Evise and Editorial Manager (EM) systems use indexing and extract terms to route manuscripts to the appropriate subject areas and editorial teams [752]. Similarly, IEEE’s Manuscript Central (built on ScholarOne) combines metadata, author-provided keywords, and an academic-network reviewer discovery tool for more accurate matching [698]. Springer’s SNAPP system and Nature’s AI-assisted triage tool also demonstrate AI’s impact on desk-review workflows [560]. Additionally, Díaz et al. [169] develop AnnotateGPT, which generates annotations to help editors quickly assess a manuscript’s scope and quality, further speeding up the review process.

7.1.2. Reviewer Matching

Reviewer matching in peer review assigns manuscripts to experts whose knowledge aligns with the submission, aiming to maximize review quality, fairness, and workload balance [592]. Charlin et al. [99] first formulate this as an integer program, using affinity scores to balance quality, fairness, and load. Charlin and Zemel [98] later embed papers and reviewer profiles in a shared latent topic space, improving efficiency and accuracy in large conferences. As submission volumes increased, automated conflict-of-interest (COI) detection became essential. Wu et al. [830] introduce a semi-automated COI declaration system and a supervised ranking model to flag conflicts and ensure fairness. Pradhan et al. [600] further advance this with a greedy algorithm that optimizes expertise distribution and workload while maintaining COI constraints, enhancing both fairness and efficiency. To address growing demands, Leyton-Brown et al. [426] develop the Large Conference Matching (LCM) algorithm, balancing expertise and load across thousands of papers. Fu et al. [214] and Aitymbetov and Zorbas [13] tackle interdisciplinary submissions by forming multidisciplinary reviewer teams to improve review quality.

7.2. In-Review

This stage involves generating or supporting review reports, either through automation or human reviewer assistance. Reviewers must assign a numerical score and provide a written evaluation. The in-review process typically involves two main stages: Peer-Review and Meta-Review.

Model	Focus similarity				Text similarity		
	KL Divergence	Overall F1	Strength F1	Weakness F1	ROUGE-L	BERTScore	BLEU-4
GPT-4o-mini [641]	0.081	0.344	0.335	0.353	0.197	0.883	0.076
GPT-4o [641]	0.082	0.348	0.342	0.354	0.202	0.885	0.079
o1-mini [334]	0.090	0.359	0.331	0.385	0.179	0.878	0.059
o1 [334]	0.097	0.355	0.318	0.388	0.170	0.869	0.032
DeepSeek-R1 [263]	0.120	0.373	0.341	0.400	0.156	0.874	0.045
Llama-3.1-70B [250]	0.136	0.339	0.338	0.341	0.215	0.882	0.076
Llama-3.1-405B [250]	0.145	0.349	0.349	0.350	0.218	0.884	0.089
DeepSeek-V3 [477]	0.151	0.350	0.330	0.368	0.199	0.880	0.069
GPT-4o-Finetuned [641]	0.022	0.306	0.280	0.322	0.194	0.882	0.081
MARG [157]	0.113	0.346	–	0.346	0.160	0.854	0.011

Table 5: Comparison of expert and LLM review performance based on Shin et al. [679], where “GPT-4o-Finetuned” refers to GPT-4o finetuned with review data using the finetune-API. KL divergences are calculated from the average of four focus distributions (strength/target, weakness/target, strength/aspect, weakness/aspect) between expert and LLM reviews. F1 scores for overall performance, strength, and weakness are derived by comparing the (target, aspect) sets between expert and LLM reviews. Text similarity metrics are computed to assess the alignment between LLM and expert reviews. Results are sourced from Shin et al. [679]. The **bolded** contents indicate the highest performance for each metric.

7.2.1. Peer-Review

Peer-review generation involves the automatic creation or assistance in the development of review comments for submitted manuscripts, including predicting quality scores and providing textual feedback. A detailed comparison of these findings is presented in Table 5.

Score Prediction estimates scores on criteria like innovation and clarity, assessing overall quality through multiple feature points [58]. Jia et al. [339] introduce a multi-task BERT framework that jointly detects quality features (e.g., suggestions, problem mentions) in review comments, outperforming single-task baselines. RelevAI-Reviewer [147] treats review tasks as a classification problem to predict papers’ relevance to a given call. Basuki and Tsuchiya [45] frame score prediction as a regression task using internal paper features, excelling at distinguishing “good” from “poor” submissions. To tackle data scarcity, Muangkammuen et al. [553] improve upon this method by introducing a semi-supervised approach that fine-tunes a transformer-based model using unlabeled data, effectively utilizing contextual cues.

Comment Generation involves generating natural-language review comments, which is the core element of manuscript evaluation [873]. Robertson [641] demonstrate that GPT-4 can generate plausible review comments. Yuan and Liu [879] construct concept graphs and integrate citation mapping on a pre-trained model to generate comments. AI-Scientist [507, 848] found that LLM-based agents approach human-level review performance [461]. MARG [157] assigns paper sections to multiple LLM agents for internal discussion, improving feedback relevance. Chamoun et al. [94] allocate four specialized roles to enhance specificity and comprehension. Furthermore, AgentReview [349] and Tan et al. [739] model the review process as a dynamic, multi-round dialogue.

Unified Generation integrates textual comments and numeric scores into a single review output that mirrors real-world peer review workflows [678, 415]. There are three main paradigms for optimizing unified peer-review generation: (1) **Single-Agent Optimization**: A straightforward approach is to optimize a single

agent through deeper analysis [358, 940, 327]. Shin et al. [679] observe that, by comparing the focus distributions of LLMs and human experts, off-the-shelf LLMs tend to prioritize technical validity in paper reviews while underemphasizing novelty. To address this, Tyser et al. [762] enhance the review system with a suite of review documents to reduce risks of misuse, score inflation, overconfident assessments, and uneven distributions. Additionally, Zhu et al. [940], Zhang and Abernethy [899] incorporate deeper reasoning via reasoning LLMs to improve review quality. **(2) Iterative Refinement Optimization:** High-quality feedback is often ensured through hierarchical quality control and multi-round refinement loops [59]. Wu et al. [828] propose an LLM-driven pipeline with hierarchical verification, producing literature surveys that match human-authored reviews. Kirtani et al. [380] introduce standardized evaluation metrics and a self-refinement cycle to align LLM-generated reviews with human accuracy and analytical depth. **(3) Multi-Agent Optimization:** To further enhance feedback reliability, some studies adopt multi-agent frameworks [739, 321, 220, 570]. D’Arcy et al. [157] divide manuscripts into modules for specialized agents, leading to higher-quality feedback than single-agent systems. CycleResearcher [820] and TreeReview [97] apply reinforcement learning to simulate iterative review rounds and structured agent interactions, enhancing collaboration. Furthermore, Taechoyotin and Acuna [735] propose multi-objective reinforcement learning to optimize unified peer review, while Taechoyotin et al. [736] extend multi-agent scientific reviews to multimodal scenarios.

7.2.2. Meta-Review

Meta-review generation synthesizes multiple reviewers’ opinions into a single, objective, and comprehensive critique, emphasizing the manuscript’s core contributions and limitations while balancing diverse viewpoints [393, 435, 296]. Early studies focus on guiding the summarization process through explicit structural cues [434, 660, 885]. More recent work addresses argumentative structures and latent biases among reviewers [113]. Notably, PeerArg [720] introduces a Multiparty Argumentation Framework (MPAF) that combines LLMs with knowledge representation to reduce subjectivity and bias. MetaWriter [725] automates the extraction of key arguments from reviewers. Darrin et al. [158] adapt the Rational Speech Act framework by creating a “distinctiveness score” to identify shared and unique perspectives across reviews. Moreover, Kumar et al. [395] introduce the ContraSciView corpus, which automatically detects contradictions between review pairs. Together, these efforts pave the way for more transparent and equitable meta-reviews.

7.3. Post-Review

Post-Review refers to the suite of AI-driven methods applied after a paper has passed peer review, aiming both to assess its future scholarly impact and to broaden its dissemination. It encompasses (1) influence analysis, predicting citation trajectories and research significance from the paper’s content; and (2) promotion enhancement, automatically generating posters, lay summaries, videos, and other outreach materials to maximize visibility.

7.3.1. Influence Analysis

Influence analysis seeks to predict the future scholarly impact of a paper, most commonly measured by citation count, by evaluating its intrinsic characteristics [694, 325]. Early approaches predominantly rely on external metadata or handcrafted features, such as author reputation and journal impact factor [172, 737, 944]. In contrast, recent methods leveraging LLMs offer the advantage of directly inferring a work’s innovativeness from its narrative. For instance, Zhao et al. [918] frame influence prediction as a regression task, fine-tuning an LLM on titles and abstracts to generate a time- and field-normalized impact score, effectively addressing the cold-start problem. Similarly, the HLM-Cite framework [275] adopts a two-stage approach: first, an embedding model retrieves a set of candidate citations from a large corpus, followed by a generative LLM that performs fine-grained reasoning and re-ranking to identify the most relevant references. Empirical

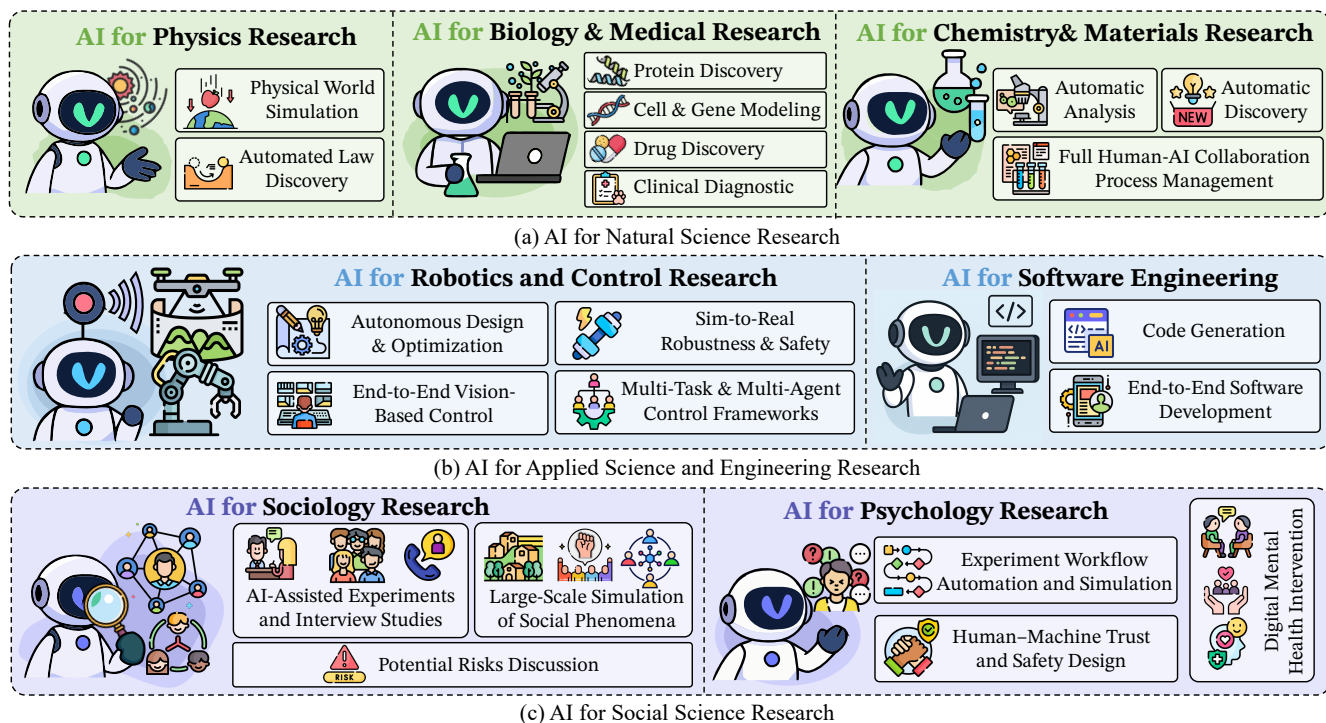


Figure 8: Multidisciplinary Applications of AI in Research. This includes three primary areas: (a) AI in Natural Sciences, covering fields such as physics, biology and medicine, and chemistry and materials science; (b) AI in Applied Sciences and Engineering, focusing on robotics and software engineering; and (c) AI in Social Sciences, encompassing disciplines such as sociology and psychology.

studies [515] suggest that these content-based methods can even surpass human experts in predicting research outcomes in fields such as neuroscience.

7.3.2. Promotion Enhancement

Beyond predicting impact, a parallel research strand employs generative AI to amplify a paper’s influence by producing varied, accessible promotional materials. These tools convert dense scientific manuscripts into more inviting formats, thereby broadening their reach. For instance, Sun et al. [728] present the P2P system, which automatically generates academic posters from lengthy, multimodal documents through intelligent content selection and optimized layout design. To improve public understanding, Markowitz [532] leverage GPT-4 to produce lay summaries and demonstrate that these AI-generated summaries surpass human-written ones in linguistic simplicity. More recently, Park et al. [590] introduce SciTalk, a multi-agent framework that generates concise scientific videos. The rapid proliferation of such systems highlights the critical need for robust evaluation, leading to the creation of specialized metrics for assessing the quality of AI-produced scientific communications [295].

8. Application of AI for Research

As shown in Figure 8, it contains three main categories: AI for Natural Science Research (§ 8.1), AI for Applied Science and Engineering Research (§ 8.2) and AI for Social Science Research (§ 8.3).

8.1. AI for Natural Science Research

8.1.1. AI for Physics Research

In physics research, AI is now indispensable for developing new methodologies and driving discoveries [175, 67]. Its applications range from automated law discovery to physical world simulation and neural operator learning, all aimed at improving simulation accuracy, speeding up computation, and revealing hidden patterns from limited data [346, 806, 542].

Physical World Simulation integrates physical priors with AI models to simulate complex systems while enforcing consistency with physical laws [47, 162, 629]. Earlier, Physics-Informed Neural Networks (PINNs) [628] embed PDE constraints in the loss function, allowing them to solve and infer nonlinear equations from sparse data. By exploiting the conserved-quantity structure of Hamiltonian mechanics, Hamiltonian Neural Networks [251] exploit conserved-quantity structures to enforce energy conservation, yielding faster convergence and drift-free, reversible simulations. Lagrangian Neural Networks [149] parameterize the system’s Lagrangian directly, avoiding coordinate choices, while still preserving exact energy conservation in examples like the double pendulum.

Automated Law Discovery leverages the reasoning power of LLMs, automated law discovery systems generate, test, and refine physical laws from noisy experimental data [681, 683]. For instance, AI-Newton [193] autonomously derives and validates physical laws, such as Newton’s laws and conservation principles, without requiring operator-provided equations. By integrating a solid knowledge base with a structured discovery workflow, AI-Newton generates interpretable models of physical phenomena. Shojaee et al. [681] propose a novel method utilizing LLMs’ scientific knowledge and code-generation capabilities to discover scientific equations directly from data. The DrSR framework [794] enhances law discovery by analyzing structural data relationships and implementing a feedback mechanism, improving performance across various domains. LLM-Feynman [706] combines automated feature engineering, LLM-based symbolic regression, and formula interpretation to extract interpretable expressions from both empirical data and domain knowledge. Recently, Li et al. [439] use enhanced visual prompting with domain expertise to uncover physical coordinates and governing equations from high-dimensional datasets more efficiently.

8.1.2. AI for Biology & Medical Research

Artificial intelligence in the life sciences and medical research uses algorithms and computational models to analyze and predict across scales [406, 444, 855, 780, 80], from molecular structures to clinical diagnostics [823, 148, 797], to accelerate drug discovery [741, 579, 259], optimize experimental workflows [732, 328, 563], improve diagnostic accuracy, and advance precision medicine [344, 364, 854, 316].

Protein Discovery. A notable example of computational innovation is Protein Discovery and protein structure prediction, which aims to predict the three-dimensional atomic structure of proteins. This is key for understanding biological functions and guiding drug design [171, 230, 232, 591]. For instance, Senior et al. [668] show that deep learning-based distance predictions significantly enhance de novo folding accuracy. The AlphaFold 2 system, developed by Jumper et al. [353], achieves atomic-level precision and has transformed structural biology since 2021. AlphaFold 3 [2] builds on this by introducing a diffusion-model architecture that predicts monomeric structures and reconstructs protein-nucleic acid and protein-ligand complexes with near-experimental accuracy. Additionally, Lin et al. [473] present a dual-task LLaMA-based framework that integrates reaction and retrosynthesis into a unified recombination-fragmentation process, generating novel compounds with strong predicted protein-binding affinities through molecular docking feedback.

Cell & Gene Modeling. A crucial area of research is cell-level modeling and gene expression analysis, aiming to simulate cellular behavior and identify activity changes under various conditions [81, 223, 647]. Several studies focus on pretraining models to improve cell or gene modeling [123, 61, 146]. However, due to the scarcity of high-quality gene and cellular data, recent work has explored data augmentation techniques to enhance AI training data [6, 528, 11]. Additionally, Roohani et al. [648] introduce an agent-based intelligent system that designs novel experiments, reasons about outcomes, and efficiently navigates hypothesis space by utilizing external tools to search the biomedical literature, analyze datasets, and engage secondary agents for evaluation, thus converging on optimal solutions. Furthermore, recent research has investigated AI-driven autonomous medical procedures, positioning AI as a collaborative tool for researchers [949, 840]. The micro-STAR system [282] integrates real-time OCT imaging with AI tissue classification to autonomously perform vascular suturing on ex vivo vessels, achieving leak-pressure performance comparable to expert surgeons, thus demonstrating the potential of AI and robotics in minimally invasive surgery.

Drug Discovery In recent years, artificial intelligence (AI) has made significant advancements in the field of drug discovery, driving multi-faceted innovations in drug design and showcasing immense potential and prospects [267, 741, 183]. **(1) Structural Prediction and Molecular Design:** AI has made notable progress in structural prediction and molecular design. Early studies [715] use deep learning models to screen 23 potential antibiotic candidates from over 107 million molecules, successfully identifying a drug with antimicrobial activity. The LUMI-lab platform [151], which integrates molecular models with automated experimentation, discovering ionized lipids that excel in mRNA delivery. However, challenges related to data scarcity persist in AI-driven drug discovery. To mitigate this, strategies such as multi-target drug polypharmacology, decoding drug responses, and quantum computing have been proposed to enhance model performance [219]. **(2) Multi-Agent Collaborative Drug Identification:** Multi-agent systems have proven highly effective in drug discovery, facilitating the rapid identification of new therapeutic compounds [417]. For instance, the DrugAgent [493] and DrugPilot [430] automate machine learning programming through multi-agent collaboration, achieving full automation from data acquisition to model evaluation, thus improving efficiency. Solovev et al. [699] introduces multi-agent approach that combines LLMs with specialized generative models and validation tools to automate the end-to-end drug discovery process. Lee et al. [417] develop a multi-agent framework that retrieves and integrates information from biomedical knowledge bases to generate responses, avoiding the need for expensive domain-specific fine-tuning. Despite these advances, challenges remain in addressing data quality, model interpretability, and regulatory hurdles [579, 878]. **(3) Drug Repurposing:** Drug repurposing involves the use of approved drugs for new therapeutic indications [509, 315]. In liver fibrosis research, the AI-assisted Collaborative Scientist system has successfully recommended drugs with significant anti-fibrotic activity using human liver organoid platforms, showing promise for treating liver fibrosis [259]. By integrating knowledge graphs and diverse data sources, Liu et al. [493] and Gharizadeh et al. [234] identify potential drug repurposing candidates and provide interpretable predictions. Lee et al. [418] combine subgroup analysis and treatment effect estimation, simulating clinical trials to identify drug candidates and characterize patient subgroups based on treatment effects. This approach, tested on a real-world database of more than 8 million patients, simulate over 1,000 drug trials, identifying 14 drug candidates beneficial to specific subgroups.

Clinical Diagnosis Clinical diagnosis advances through three converging breakthroughs. **(1) Clinical Brains:** First, LLMs serve as clinical “brains”, matching physician-level performance on medical licensing examinations [77, 249, 372, 692, 823]. Further, to enhance transparency and structure in decision-making, a human-AI note-taking framework [148] and Long-CoT reasoning techniques [409] has been proposed, leveraging case-based reasoning to guide clinical inquiry. **(2) Multi-Agent Hospital Simulation:** Multi-agent systems such as Agent Hospital replicate AI-AI [192, 401] and human-AI [662] collaborative diagnostic and treatment workflows, effectively serving as an organizational “nervous system” for care coordination [429, 39].

(3) Interactive Physical Actuation: Robotic platforms guided by LLMs perform precise physical interventions. For example, an autonomous optical coherence tomography system delivers surgeon-level accuracy in delicate procedures such as vascular anastomosis [282]. Together, these breakthroughs demonstrate the feasibility of fully autonomous medical facilities in which artificial agents seamlessly integrate diagnostic reasoning, therapeutic planning, and procedural execution.

8.1.3. AI for Chemistry & Materials Research

AI-driven automation in chemistry [283, 518, 552] and materials [610, 345, 134] integrates machine learning, robotics, and instrumentation into a closed-loop system for design, synthesis, and characterization, speeding decisions and experiments [115, 881, 827].

Automatic Analysis seeks to identify optimal or novel material compositions in virtual or automated experimental setups while minimizing the number of required experiments [46, 630, 669]. Specifically, Chen et al. [101] introduce MEGNet, demonstrating that graph neural networks can achieve density functional theory-level accuracy for both molecular and crystalline properties. Li et al. [450] employ two-stage Bayesian optimization to screen 560 organic photocatalysts using only 2.4% of the experimental conditions, thereby obtaining significantly improved performance. Ekosso et al. [186] combine low-cost robotic platforms with high-throughput microscopy and Gaussian process models to map vesicle formation processes. More recently, Szymanski et al. [734] implement a robot-machine-learning platform that accelerate compound discovery and identify 41 novel inorganic materials within 17 days.

Automatic Discovery is an automated experimental platform that combines robotic operations, online characterization, and real-time decision-making algorithms to autonomously execute the full experimental process, from reagent dispensing to result analysis [83, 458, 159]. Early research laid both theoretical and practical foundations: Butler et al. [83] review machine learning methods that accelerate materials design and discovery [546, 156, 322]. Furthermore, Dai et al. [153] employ a mobile robot with UPLC-MS and NMR to plan and interpret syntheses in a manner akin to human chemists. Jayarathna et al. [337] leverage literature data to reduce the number of experiments in an active-learning loop, discovering new Ru-based catalysts. Dai et al. [154] introduce an AI advisor for ion-electron polymers, improving performance by 150% over spin-coating methods. More recently, several studies have incorporated LLMs to enhance innovation and knowledge in chemistry and materials discovery [902, 857, 388].

Full Human-AI Collaboration Process Management leverages LLMs or natural-language understanding and generation to support hypothesis formulation, experimental design, and iterative optimization, aiming to facilitate more intuitive and efficient research interactions [781, 513, 422]. The AILA framework [529] embeds LLMs within a fully automated atomic-force microscopy workflow, illustrating both the potential and current constraints of language models in guiding real-time microscopic experiments. Sprueill et al. [708] and Feng et al. [203] integrate LLM-driven linguistic reasoning with chemical feedback in a heuristic search loop to propose novel catalysts and reaction pathways within uncertain chemical spaces. Recognizing that collective intelligence often outperforms individual reasoning, several studies [566, 704, 404] employ multi-agent architectures in which LLMs generate hypotheses, design experiments and direct iterative optimization, thereby achieving seamless human-AI collaboration. Ma et al. [520] introduce the first fully automated retrosynthetic planning agent tailored for LLM-driven macromolecular design, enabling comprehensive enumeration of viable multi-branch reaction routes. Meanwhile, Zhu et al. [942] demonstrate a robotic AI chemist that performs ore pretreatment and catalyst optimization on Martian meteorite samples.

8.2. AI for Applied Science and Engineering Research

8.2.1. AI for Robotics and Control Research

AI for Robotics and Control Research applies AI methods: deep learning, reinforcement learning, and large language models, to the perception, decision-making, and control of robots, aiming to boost adaptability, robustness, and autonomy in novel environments [419, 411].

Autonomous Design & Optimization systems integrate robotics, machine learning, and domain expertise to automate experiment planning, execution, and optimization [458]. Uddin et al. [763] introduce OptoMate, a system using a fine-tuned language model for optical setup design and a robotic arm to assemble spectroscopy components with submillimeter precision, enabling cloud-based optical labs. Mieszczanek et al. [549] employ computer vision and feedforward neural networks in a feedback optimizer to adjust 3D printing parameters in real-time, reducing data collection from days to hours and ensuring consistent part quality. Angello et al. [22] apply physically informed feature selection and supervised learning in a closed-loop system to enhance photostability and uncover solvent-mediated triplet-state mechanisms. Bu et al. [78] combine a text-conditioned video diffusion model with a feedback-driven controller to generate visual plans and iteratively refine actions, significantly improving performance.

End-to-End Vision-Based Control End-to-end vision-based control feeds raw images or video frames directly into a neural network to generate control signals, removing separate perception, planning, and control modules. Early work combine convolutional neural networks (CNNs) with reinforcement learning or guided policy search to map camera inputs to motion commands. Levine et al. [424] first apply Guided Policy Search to train perception and control jointly, mapping raw images to motor torques and demonstrating reliable real-world grasping. Levine et al. [425] train a CNN on massive real grasp attempts to predict grasp success in real time, closing the loop on novel objects. Kalashnikov et al. [356] introduce a self-supervised, closed-loop Q-learning framework trained on 580,000 grasps, enabling dynamic strategy adjustment, retries, and disturbance resilience. Tobin et al. [755] propose domain randomization, varying simulator rendering parameters so models trained on synthetic data transfer directly to real-world detection and grasping.

Sim-to-Real Robustness & Safety ensures reliable transfer of simulation-trained policies to real-world tasks while adhering to safety constraints. Bochem et al. [66] integrate sharpness-aware optimization into gradient-based RL to identify flat minima, enhancing transfer robustness in contact-rich tasks without compromising sample efficiency. Ayabe et al. [34] assess offline RL methods on a legged robot subjected to random and adversarial torque disturbances, revealing vulnerability to sudden perturbations and emphasizing the need for real-time adaptation and safety measures. Radosavovic et al. [624] train a Transformer-based controller using deep RL and deploy it outdoors for one week without safety scaffolding, showcasing adaptive performance amidst disturbances, rugged terrain, and varying payloads. Yang et al. [852] apply domain randomization for vision-based servoing of soft robots, eliminating the need for on-robot fine-tuning and enabling direct transfer of simulation-trained models to continuum manipulators. Guerrier et al. [261] combine control barrier functions with RL to enforce safety constraints during learning, preventing hazardous states in complex environments.

Multi-Task & Multi-Agent Control Frameworks facilitate concurrent task execution or enable collaboration among agents in complex workflows, thereby enhancing parallelism and automation. Tahmid and Notomista [738] introduce a reinforcement learning-based framework designed to dynamically learn and compose task policies in robotic systems with redundant architectures, incorporating time-varying priority stacks to adjust task priorities. Team et al. [749] propose a unified multi-agent system capable of automatically generating hypotheses, designing and conducting experiments, and refining methods through iterative feedback, establishing a closed-loop process that accelerates interdisciplinary research.

8.2.2. AI for Software Engineering

AI for Software Engineering Research focuses on applying AI techniques to automate software development tasks, enhance code quality, and improve developer productivity. This includes code generation, bug detection, code review, and software testing.

Code Generation refers to the use of AI models to automatically generate code snippets or entire programs from natural language descriptions or existing code patterns [654, 262, 436, 313]. This can accelerate the development process and reduce manual coding [506, 912]. For instance, Chen et al. [105] develop Codex, a GPT model fine-tuned on GitHub’s publicly available code, which supports GitHub Copilot. To democratize program synthesis, Nijkamp et al. [569] train and release CodeGen, an LLM based on both natural and programming language data, alongside the open-source training library JAXFORMER. Additionally, several studies have explored advanced code capabilities and support for multiple programming languages, like Python, Java, and R [654, 262, 436].

End-to-End Software Development covers the entire software development lifecycle, with AI automating various stages [581, 190, 347, 413]. For example, Phan et al. [596] develop HyperAgent, a generalist multi-agent system designed to handle various SE tasks across different programming languages, mimicking human developers’ workflows. Qian et al. [611] introduce Experiential Co-Learning, which enables software development agents to leverage historical experiences to improve task performance. Meanwhile, Qian et al. [612] introduce ChatDev, a chat-based framework for software development, and Kang et al. [361] present an explainable automated debugging framework powered by LLM-driven scientific debugging.

8.3. AI for Social Science Research

AI has been widely utilized to automate the design, execution, and analysis of social science experiments, encompassing tasks from hypothesis generation to data collection, with minimal human intervention. In this context, we will focus on two key domains:

8.3.1. AI for Sociology Research

AI in sociology research refers to the use of machine learning, natural language processing, and multi-agent systems to simulate, analyze, and explore social phenomena [504, 838, 365]. Through AI, researchers can reconstruct macro-level patterns of collective behavior and gain deeper insights into micro-level cultural contexts and individual interactions, thereby revitalizing traditional sociological methods [456].

AI-Assisted Experimental and Interview Studies. Controlled experiments and simulated interviews are increasingly employed by scholars to test social science hypotheses and evaluate the effects of various social mechanisms and policy interventions. Manning et al. [530] propose a methodology that combines structural causal models with large language models to automatically generate and empirically validate social science hypotheses in contexts such as negotiations, bail hearings, job interviews, and auctions. This approach effectively bridges the gap between theory and practice by utilizing the model both as a scientific tool for hypothesis generation and as an experimental subject for validation. Liu and Yu [482] develop MimiTalk, an automated interview system, and conduct a comparative study of AI-led and human-led interviews with 20 participants on the Prolific platform. This study demonstrates the feasibility of AI-mediated interviews and highlights their potential in experimental settings.

Large-Scale Simulation of Social Phenomena. This approach leverages algorithmic tools to automate the collection and analysis of extensive textual, visual, and interaction data to simulate and examine the macro-level dynamics of community practices and value evolution [292, 710, 324]. Perez et al. [595] automate

the extraction and analysis of large-scale text and image datasets to map cultural practices, value systems, and trends in contemporary online communities [871, 674]. Zamudio et al. [883] propose a simulation framework for cultural evolution using multi-agent LLMs, enabling the manipulation of network structures, individual traits, and biases in information transmission to investigate factors driving cultural diffusion and change. Chen et al. [120] develop a GPT-based three-module framework, including information extraction, variant generation, and outcome prediction, that achieved high consistency in predicting outcomes across 319 economic field experiments, while also reflecting the impact of gender, race, and social norms on performance. Additionally, Bao et al. [38] reveal the underlying, often unspoken codes within societies.

Potential Risks Discussion. While LLMs demonstrate strong predictive capabilities in the natural sciences, their performance in the social sciences remains limited. Manning et al. [530] find that, although LLMs can predict the signs of estimated effects well when given a proposed structural causal model, they struggle to predict the magnitudes reliably. Additionally, Luke et al. [510] highlight that LLMs face challenges in handling treatment effect heterogeneity and exhibit systematic biases when predicting social science outcomes. As such, LLMs’ predictive capabilities are still underdeveloped, particularly in forecasting novel empirical patterns that could inform future experimentation [421].

8.3.2. AI for Psychology Research

Research methodology focuses on the design, implementation, and validation of psychological experiments to ensure validity and reproducibility [757, 427].

Experiment Workflow Automation and Simulation. Recent research has explored integrating AI into the management and data simulation of psychology experiments [619, 784, 64]. Zamudio et al. [883] introduce the RAISE pipeline, automating the generation and validation of visual stimuli. In five experiments, AI-generated images match researcher-designed stimuli in both validity and recognizability. Cingillioglu et al. [139] conduct a fully automated online RCT with 1,193 participants, where AI managed recruitment, random assignment, intervention delivery, and data collection, successfully replicating eight classical hypotheses with gold-standard rigor. Cui et al. [152], Strachan et al. [716], Suri et al. [730] use GPT-4 to simulate responses for 154 classical experiments, reproducing 76% of primary effects but yielding 71.6% unexpected significant outcomes, illustrating the promise of AI-assisted replication while emphasizing the need for cautious interpretation [124, 170, 242].

Human-AI Trust and Safety Design. Research on human-AI trust explores the development of trust during human-AI interactions and derives strategies for ensuring safety [614]. Li et al. [452] introduce a three-dimensional framework encompassing the trustor, trustee, and context, identifying key factors influencing trust and offering design recommendations for improving safety. Building on this, Chandra et al. [96], through interviews with 283 individuals who have mental health experiences, develop a taxonomy comprising 19 risky AI behaviors and 21 negative psychological impacts. From this, they propose a multi-path case-method framework and a set of safety guidelines aimed at mitigating these risks.

Psychological Interventions. Psychological interventions increasingly employ AI-driven chatbots to provide scalable and cost-effective psychological support [604]. Earlier, Hagendorff et al. [271], Dillion et al. [170], and Binz and Schulz [64] discuss whether and when LLMs can replace human participants in psychological research, reviewing early evidence and proposing a theoretical framework while highlighting methodological caveats. In a randomized controlled trial, Heinz et al. [287] find that the Therabot chatbot resulted in significant reductions in clinical-level symptoms compared to the control group. Similarly, Spyska [709] explore the use of the Friend chatbot for crisis support, demonstrating that its efficacy is comparable to traditional face-to-face therapy. These findings highlight the potential of generative AI to enhance accessibility

Tool	Description
SciSpace Copilot	AI-powered Literature Q&A, Annotations, Auto-Summarization, Chart Explanations
Elicit	AI-powered Literature Q&A, Auto-Summarization, Suggestions
Jenni AI / NoteGPT	AI-powered Note-Taking, Auto-Summarization
Scholarcy	AI-powered Auto-Summarization, Summarization Card, Analysis and Organization
PDFMathTranslate	AI-powered PDF Math-Augmented Translation

Table 6: Representative and established AI systems and assistant tools for advancing scientific comprehension.

and effectiveness in mental health services, particularly in settings with limited resources [164, 165]. Recent work further demonstrates that LLMs can match or exceed human performance in generating emotionally resonant narratives [819, 651] and even pass standard Turing tests [352], underscoring their broader psychological and communicative capabilities.

9. Resources

To further advance research in this field, we will provide an expanded and more comprehensive suite of relevant resources, including tools, benchmarks, and datasets spanning all stages.

9.1. AI for Scientific Comprehension

9.1.1. Textual Scientific Comprehension

To advance the evaluation of scientific question-answering systems, various benchmarks have been developed with increasing task complexity and domain specificity [650, 856, 201, 285]. Table 6 presents a comprehensive overview of typical, mature AI systems and associated tools for scientific comprehension.

Datasets like ScienceQA [659], LitQA [405], LitQA2 [696], SciQA [32], SciQAG-24D [771], and TriviaQA [391] support QA for scientific content. SciBench [798] broadens scientific reasoning across physics, chemistry, and mathematics. Further, SciInstruct [889] broadens reasoning across formal proofs with instruction-tuned data [811, 769]. Moreover, AutoPaperBench [377] and SciCUEval [874] are proposed for automatic paper or scientific content understanding evaluation.

Furthermore, datasets have expanded into broader domains, including biomedicine [71, 348, 389, 636, 584, 472, 915], academic chemistry [119, 594], materials science [514], physics [928] and other scientific fields [845]. TheoremQA [114] evaluates AI models’ ability to apply theorems to solve challenging science problems. Multi-task and multi-modal assessment frameworks, such as M3CoT [108], SciFIBench [640], MMSCI [453], SPIQA [602], and MultimodalArxiv [431], further extend these evaluations. To address broader multimodal and multi-document challenges, M3SciQA [428], SceMQA [463] and SciDQA [691] have also been introduced.

Beyond static benchmarks, dynamic and interactive evaluation frameworks have emerged. SCITOOL-BENCH [522] target tool use in scientific reasoning across domains, while Kuhn et al. [391] propose a multi-round dialog framework to simulate user interactions, introducing metrics like adjusted accuracy. In terms of generation alignment, Yu et al. [869] design a system to evaluate the semantic fidelity between generated content and scientific texts, combining automated scores with human judgment.

Tool	Description
Google Scholar / Web of Science / Scopus / AMiner	Literature Search, Citation Tracking, Author Profiles, Citation Analysis
Semantic Scholar	AI-Assisted Academic Search Platform (Semantic Graph)
Research Rabbit / Connected Papers / Citation Gecko / Iris.ai	Visual graph of Works
Scite.ai	Shows Citation Context (Supporting/Contradicting/Neutral)
Consensus.app	Opinion-based Literature Search, Ideal for YES/NO Questions
ResearchGPT	AI-generated Knowledge Graphs and Paper Structures

Table 7: Representative AI systems and assistive technologies that have been widely adopted to enhance academic surveys.

9.1.2. Table & Chart Scientific Comprehension

In the domain of reasoning based on charts and tables, a number of benchmarks have emerged to evaluate the ability of LLMs in both structural and logical comprehension. Early works, such as ChartQA [535], CharXiv [812], ChartX [837], and NovaChart [301], focus on assessing LLMs’ performance in answering questions related to charts, utilizing both synthetic and real-world data. On the other hand, benchmarks like SUC [718], TableBench [832], and ToRR [28] emphasize the evaluation of LLMs’ structural understanding, rather than content comprehension, across various tasks such as form interpretation, numerical reasoning, and textual analysis.

9.2. AI for Academic Survey

In the task of generating academic surveys, several representative public corpora stand out due to their distinct characteristics in terms of scale, domain, and structure. To evaluate the capabilities of scholarly retrieval systems, several studies have focused on the scholarly deep research [815], such as AcademicBrowse [932]. To facilitate section-level generation of related works, several large-scale datasets have been introduced, such as Cochrane [770], MSLR 2022 [787], MS² [168], and OARelatedWork [173], OAG-Bench [891]. These datasets pair comprehensive “Related Work” sections with their corresponding full texts, providing valuable resources for this task. Moreover, systematic benchmarks for evaluating the quality of automatic scholarly survey generation have been developed. Examples include SciReviewGen [366], BigSurvey [491], SurveySum [205], SurveyBench [849], AutoSurvey [802], and SurveyX [462]. These benchmarks provide critical metrics for assessing the performance of automatic systems in generating academic surveys. For fine-grained manual annotation, SurveyEval [782] offers a hierarchical title tree that includes a vast number of reviews and citations. It is accompanied by hierarchical consistency and citation-chapter alignment metrics, which serve as essential tools for evaluating the distribution of synopsis generation and citation accuracy.

9.3. AI for Scientific Discovery

Idea Mining has seen the introduction of several key resources that significantly contribute to scientific discovery tasks [104]. Notably, LiveIdeaBench [655], ResearchBench [500], Genome-Bench [868], AllIdeaBench2025 [620], the AP-FRI Corpus [396], HypoGen [575], CLIMATEDATABANK [495], CHIMERA [714] and OMATO-Chem [861] provide structured datasets for idea mining and hypothesis generation, enabling systematic training and evaluation of LLMs. Furthermore, OAG-Bench [891] offers a comprehensive, fine-grained benchmark for academic graph mining, spanning 10 tasks, 20 datasets, and over 70 baseline methods, all curated by human experts. This resource fosters systematic evaluation and encourages community-driven research. Additionally, SPARK [661] and the ICLR-NeurIPS Ideas Dataset [440] introduce curated datasets of idea-centered abstract-review pairs from OpenReview submissions, supporting supervised and reinforcement learning for research idea generation with multi-dimensional idea quality control, including novelty,

Tool	Description
<i>Experiment Design</i>	
SnapGene Elicit	Molecular Cloning and DNA Visualization AI-driven Experiments Design
<i>Experiment Management</i>	
Notion / Asana / ClickUp / Atlassian Rovo Trello / Wrike GitMind Forecast Tableau / Power BI	Project Management, Task Tracking, Collaboration Generate Content, Help Brainstorming, Conceive Product Mind Mapping and Brainstorming Tool for Project Planning Project Risk and Status Management Interactive Dashboards and Reports
<i>Experiment Conduction</i>	
Copilot / Cursor / Tabnine / Qodo Gemini CLI Diffblue MLflow / Weights & Biases / TensorBoard Papers with Code MONAI Taskade	AI-powered Code Completion, Generation, Review, Documentation AI-powered Open Source Command Line Tool AI-powered Unit Test Generation for Java Code AI Experiment Tracking And Visualization Paper-Code Pairs for Easier Reproducibility AI-powered Medical Imaging Framework for Reproducibility Generate Code Snippets and Debugging to Facilitate Collaboration Among Developers
<i>Full-Automatic Discovery</i>	
ChatGPT / Claude / Gemini ResearchGPT AutoGPT / OpenDevin AgentLabs AI-Scientist / Zochi	Problem Solving, Code Assistance, Writing Polishing, Full-Lifecycle Management AI-generated Knowledge Graphs and Paper Structures Multi-step Research Automation Multi-Agent AI Platform for Research Automation AI-powered Research Assistant for Scientific Discovery

Table 8: Representative AI tools and their role in facilitating scientific discovery, with a particular focus on experimental conduction.

feasibility, and effectiveness.

Novelty & Significant Assesment In the development of automated scientific research evaluation, the academic community has primarily focused on the dual criteria of “novelty and significance”, systematically exploring the ability of language models to assess the innovation within scientific research [740]. The SchNovel framework [467] and NoveltyDetection [497] are introduced to evaluate AI systems’ capacity to assess scholarly novelty across multiple scientific disciplines sampled from arXiv, aiming to facilitate the automated evaluation of research originality in scientific workflows. Building upon this, Gu et al. [254] introduce BLADE, a system that integrates 12 expert-labeled datasets with multiple automated scoring methods, enabling the model to explore diverse inference strategies in open, data-driven scientific analysis. HypoBench [484], Dasgupta et al. [161] and Lin et al. [471] are designed to evaluate LLMs and hypothesis generation methods across multiple aspects, including practical utility, generalizability, hypothesis, novelty [161] and rigor [471] discovery rate.

Theory Analysis requires the collection of scientific evidence, theoretical verification, and theorem proving. Specifically, FV-Generalization Benchmark [586], SCitance [17], and MissciPlus [240] are designed to complete scientific evidence collection. TheoremExplainBench [390], XClaimCheck [363], SciNews [91], ClaimReview2024+ [75], FactKG [376], TrendFact [900] provide datasets and benchmarks for scientific verification analysis. MiniF2F [925], FIMO [480], MUSTARDSAUCE [323] are used to fine-tune and evaluate LLMs on scientific theorem proving tasks. Furthermore, datasets and benchmarks have expanded into broader domains, including biomedicine [809].

Tool	Description
EndNote / Mendeley plugins	Reference Insertion and Auto-Formatting
Mathpix Snip / MathHandwriting	AI-powered Math Equation Recognition and LaTeX Conversion
AI for Grant Writing	AI-powered Grant Writing Assistance
Writefull / Trinko / Grammarly AI / Paperpal / Overleaf Copilot / Wordtune	AI-powered Scientific English Polishing Tools
SciSpace Copilot / Jenni AI	AI Writing Assistants for Editing and Suggestions
ChatGPT / Claude / Gemini	Writing Inspiration, Summarization, Editing
GPT-4o / Vizcom / Illustrae / OpenArt	AI-powered Figure Generation and Illustration Tools

Table 9: An overview of representative AI tools and their contributions to enhancing academic writing.

Experiment Design In terms of experimental design, Tian et al. [754] propose evaluation frameworks for zero-shot and few-shot scenarios within virtual screening and lead compound optimization, establishing a comprehensive set of metrics tailored for AI-driven drug discovery. Concurrently, Feng et al. [199] leverage 1.6 million bioactivity measurements to train a universal model using pairwise meta-learning, which facilitates rapid adaptation and robust generalization to new biological systems. For biological protocol understanding and reasoning, BioProBench [501] is the first large-scale, multi-task benchmark. In order to provide valuable insights for the safe and effective deployment of LLMs in medical domains, Zhang et al. [896] develop LLMEval-Med, a real-world clinical benchmark for medical LLMs with physician validation.

Experiment Conduction To evaluate model performance in realistic research environments, MAgent-Bench [320], Exp-Bench [383], MLRC-Bench [911], MLE-Bench [95], DS-Bench [350], ScienceBoard [727], AutoReproduce [921], SciReplicate-Bench [839], DO Challenge [697], and MLR-Bench [102] assess AI agents’ abilities to perform typical research tasks, such as optimizing CIFAR-10 classifiers and tuning BabyLM. In a similar vein, Hu et al. [307] develop InfiAgent-DABench, which is based on real-world CSV datasets and evaluates models’ ability to interact with tools in end-to-end data analysis tasks. MLGym-Bench [559] is the first Gym environment for machine learning tasks, enabling research on reinforcement learning algorithms for training such agents. ResearchCodeBench [310] enables continuous understanding and advancement of LLM-driven innovation in research code generation. AutoBio [410] is designed to evaluate robotic automation in biology laboratory environments.

Experimental Analysis Experimental Analysis involves systematically testing hypotheses, evaluating models, or validating theoretical assumptions to draw meaningful conclusions. MicroVQA [82] is proposed to assess three reasoning capabilities vital in research workflows: expert image understanding, hypothesis generation, and experiment proposal.

Full Automatic Discovery In recent years, benchmark suites have been developed to assess AI-driven research agents. These suites offer standardized datasets, predefined tasks, and evaluation metrics, thereby facilitating systematic advances in algorithmic optimization. They encompass multi-domain scenarios spanning chemical synthesis, materials discovery, and biological experimentation [254, 264, 479]. Notable examples include ScienceAgentBench [126], BaisBench [511], Curie [382], which are designed to evaluate AI scientists’ abilities to generate novel discoveries in different disciplines through data analysis and reasoning with external knowledge [697, 682, 362]. And DiscoveryWorld [336] evaluates end-to-end scientific discovery agents, while DiscoveryBench [527] challenges large language models with 264 real-world and 903 synthetic tasks across six domains, using structured protocols to measure multi-step, data-driven discovery and to elucidate both capabilities and failure modes.

9.4. AI for Academic Writing

The field of AI in academic writing is supported by a comprehensive array of meticulously curated datasets that address various aspects of the academic writing process.

9.4.1. *Semi-Automatic Academic Writing*

Recent advancements in semi-automatic academic writing have led to the development of several datasets designed to assist researchers in different stages of manuscript preparation, writing, and editing.

Assistance During Manuscript Preparation. In the early stages of manuscript preparation, recent datasets such as MoDeST [84] and LLM-Rubric [279] offer valuable tools for generating multi-domain scientific titles and assessing the scientific idea generation capabilities of LLMs.

Assistance During Manuscript Writing Several datasets, including FigGen [643], Figuring out Figures [90], SciCapenter [300], and TikZero [55], support figure and formula generation, from text-to-figure creation to automated TikZ code generation. For citation management, datasets like CITEWORTH [822], CiteBART [92], and ScholarCopilot [808] enhance context-aware automatic citation generation [106]. Additionally, FutureGen [35] extracts future work statements from thousands of papers, using LLMs to identify and validate forward-looking scientific content.

Assistance After Manuscript Completion. Once a manuscript is completed, further enhancement can be achieved through grammar correction and expression optimization. To support this process, datasets from the Automated Writing Evaluation (AWE) system [783], and AAAR-1 [505] provide valuable resources. Additionally, the transformer-based Feedback Dataset [930] offers comprehensive support for multi-dimensional writing quality assessment [333, 196]. Moreover, datasets such as Wikipedia Revision Histories [73], which track real-world editing histories, play an important role in refining language and improving overall clarity. Pang et al. [589] introduce the first benchmark and metric suite for poster generation for visual quality, coherence-language fluency, and the ability to convey core paper content.

9.5. AI for Academic Peer Reviewing

Research on AI for Academic Peer Reviewing is grounded in diverse datasets, addressing tasks from AI text detection to review generation, quality assessment, and decision support [934, 194, 632, 137, 136, 712, 48]. To simulate realistic peer review interactions, datasets such as PeerRead [358], SPOT [700], NLPeer [180], ReviewMT [739], MOPRD [470], OpenReviewer [327], MASSW [901], COMPARE [690], PeerArg [720], Re² [890], ReviewEval [380], AAAR-1 [505], Papereval [321], ORB [733] and ORSUM [886] collect extensive paper and review data from leading conferences or journals, enabling the training and evaluation of LLMs in multi-turn, long-context, or role-based peer reviews [598]. Additionally, LLMart [461] offers a toolkit for evaluating LLM robustness through adversarial testing and prompt optimization, ensuring AI reliability in sensitive academic contexts. To assess LLM review quality more precisely, Shin et al. [679] and Couto et al. [147] analyze the quality of peer review content across multiple predefined aspects, highlighting discrepancies between LLM and human review focus.

In the field of review quality detection, researchers have investigated diverse quality features [235]. Purkayastha et al. [608], and PolitePEER [60] assess AI systems' ability to identify instances of "lazy thinking" or politeness in peer reviews. Furthermore, both the AI-Peer-Review-Detection-Benchmark [876] and TRIED [490] include thousands of AI-generated peer reviews alongside human-authored reviews from the ICLR and NeurIPS conferences. These datasets provide standard corpora essential for evaluating methods designed to detect AI-generated peer reviews.

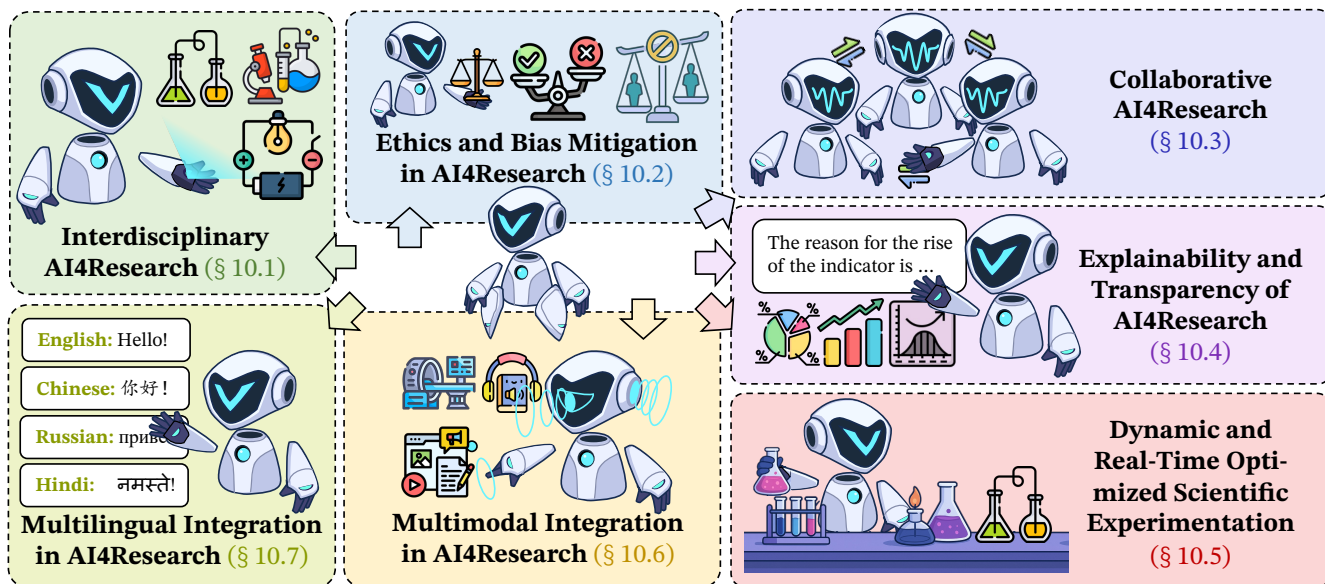


Figure 9: Frontiers and Future Directions of Artificial Intelligence in Research: This includes (1) Interdisciplinary AI models, (2) Ethics and Safety in AI4Research, (3) AI for Collaborative Research, (4) Explainability and Transparency of AI4Research, (5) Dynamic and Real-Time Optimized Scientific Experimentation, (6) Multimodal Integration in AI4Research, and (7) Multilingual Integration in AI4Research.

10. Frontiers & Future Direction

10.1. Interdisciplinary AI Models

As AI advances across research domains, we need models that integrate knowledge from multiple fields. Future work should develop general-purpose AI systems able to understand and generate insights in biology, physics, social sciences, and beyond. The primary research directions are: **(1) Foundation Models.** This paradigm has become the cornerstone of cross-domain AI. These models are pretrained self-supervised on vast unlabeled or weakly labeled datasets, then fine-tuned on new tasks with minimal data. They have driven performance gains in medical imaging, natural language processing, and robotics [312, 373]. **(2) Graph Models.** Graph methods naturally handle relational data by propagating information along nodes and edges. This enables cross-field knowledge flow, e.g., integrating ontologies and neural graphs in medical text classification for precise concept capture and efficient inference [189, 79, 257, 472, 408].

The greatest challenges at present are: **(1) Heterogeneous Interdisciplinary Data.** Interdisciplinary research involves diverse modalities, from high-dimensional sensor signals to categorical labels and unstructured text. These sources vary in scale, noise characteristics, and missing-data patterns, hindering unified preprocessing and feature fusion [609, 864, 141]. **(2) Cross-Domain Knowledge Transfer.** Transferring knowledge across domains requires extracting and adapting relevant information for new tasks. Techniques such as policy transfer, domain-adversarial training, and semantic alignment can narrow some gaps, but negative transfer persists in highly heterogeneous settings [572, 670]. Moreover, preserving reliability and interpretability during transfer, to ensure more applicable and trustworthy application in novel contexts, remains an urgent open problem [943].

10.2. Ethics and Safety in AI4Research

As AI assumes a central role in scientific research, a range of ethical, safety, fairness, and bias concerns has emerged [206, 865, 238, 317], making mitigation essential [18, 246, 843, 587, 130]. Early work by Farber [194] shows that, while AI improves reviewer matching and response rates, it disadvantages authors in low-resource languages or on niche topics. Worse, text-similarity matching can be abused by collusive rings to manipulate peer review, underscoring the need for built-in anti-collusion measures [626, 298]. Moreover, McShane et al. [539] find that AI-assisted statistical interpreters fall prey to “dichotomous mania”, reducing results to simply significant or not, a flaw that prompt-engineering alone has not resolved [634]. There are two main mitigation strategies: **(1) Fairness-Aware Training:** Integrate fairness constraints into the loss function to balance accuracy and equity across groups [206, 274]. Causal-inference methods then detect and adjust for hidden biases, enabling counterfactual fairness interventions [128, 545]. **(2) Training-free Debiasing:** Without retraining, apply unsupervised pruning and reweighting to model outputs at regular intervals, correcting biases in large language models by leveraging their pretrained behaviors [181, 649, 129]. **(3) Establishing Ethical Framework:** Some studies are establishing benchmarks for professional and broad ethical frameworks to regulate AI-generated content in a controlled area through security risk and ethical issue monitoring [939, 707].

Nonetheless, these endeavors confront two core challenges: **(1) Balancing performance and fairness:** The inherent tension between maximizing predictive accuracy and enforcing fairness constraints is difficult to reconcile and typically demands meticulous, application-specific tuning to avoid degrading model utility [289, 723]. **(2) Avoiding AI Plagiarism:** A major ethical concern in AI-driven scientific research is plagiarism [475, 217, 821]. Large-scale text generation by LLMs could lead to a “plagiarism singularity”, where text originality is diminished, raising concerns about the ethical and copyright risks of AI-generated content [631]. Studies have also revealed significant instances of intelligent plagiarism in LLM-generated scientific literature [268].

10.3. AI for Collaborative Research

As interdisciplinary research advances, the diversity of team members’ backgrounds can impede information flow and decision coordination. AI techniques can automatically extract and synchronize cross-document and cross-domain information, thereby narrowing the information gap among collaborators [65, 664]. Simultaneously, AI-driven arbitrators within real-time collaboration platforms can adjust task allocation dynamically based on project progress and member expertise, improving both efficiency and the quality of innovative outcomes [507, 198, 297]. The main research directions can be broadly divided into two categories: **(1) Collaborative Agents and Cooperative Intelligent Systems.** Collaborative agents are AI systems endowed with decision-making, autonomous execution, and communication capabilities. They simulate and augment human collaborators by participating in complex project management and research workflows through task assignment and autonomous role switching [297, 304]. Through semantic retrieval, reasoning validation, and context awareness, multi-agent frameworks are creating a new paradigm of human-AI collective intelligence, enabling automated hypothesis generation, experimental design planning, and preliminary results analysis to accelerate scientific discovery [948, 430, 710]. These advances support efficient human-AI hybrid teams and suggest fertile directions for further work on collaborative agents and distributed modeling. **(2) Federated Learning and Distributed Modeling Mechanisms.** Because sensitive data across institutions cannot be fully shared, recent research has adopted federated learning as a privacy-preserving distributed modeling approach. By training models collectively while keeping data local, federated learning mitigates data silos among institutions and specialist teams [432, 888, 398]. To enhance both performance and privacy guarantees, differential privacy, and homomorphic encryption are being integrated with federated optimization algorithms, offering scalability and regulatory compliance for

large-scale, multi-scenario collaborative research [760].

Current challenges in this field include: **(1) Interaction Complexity.** Repeated task reassignments, control handovers, and heterogeneous communication modalities can lead to misunderstandings, inefficiencies, and compounded coordination errors [243, 294]. Addressing this issue requires adaptive collaboration mechanisms that allow AI systems to adjust their behavior dynamically to match human collaborators' working styles and decision-making preferences. Multi-intelligence relationships are also critical, with three failure modes of miscoordination, conflict, and collusion [272]. **(2) Tension between Data Privacy and Accessibility.** A fundamental tension exists between data privacy and accessibility: stringent anonymization or legal restrictions often reduce the quality and diversity of training data. Although anonymization techniques and compliance with regulations protect privacy, they can diminish data utility and hinder AI models from capturing representative features, thereby affecting the accuracy and credibility of interdisciplinary research [555]. Moreover, differences in data access permissions, network bandwidth, and legal frameworks across institutions can cause communication delays and inconsistent model updates during distributed training, undermining the efficiency and stability of federated learning [260].

10.4. Explainability and Transparency of AI4Research

As AI models increasingly drive scientific discovery, ensuring their transparency and explainability is essential. Future work should strengthen model interpretability so that researchers can trace how conclusions and recommendations are generated, particularly in high-stakes scientific applications [211, 185]. Efforts to improve explainability fall into two main categories: **(1) White-box Analysis:** This approach investigates the model's internal structure by linking specific network "circuits" to conceptual representations. It has attracted considerable interest from both the security and transparency communities [57, 916, 627]. **(2) Black-box Analysis:** More recent work focuses on interpreting models without direct access to internal parameters. By examining reasoning trajectories and aggregate behavior, black-box methods provide insights into a model's knowledge representation and enable more reliable control over its outputs [85, 280, 107, 109, 111].

Despite these advances, two principal challenges remain: **(1) Lack of Standardized Frameworks:** Explanation techniques and metrics vary widely across the AI4Research community. Such absence can produce conflicting results and undermine user confidence. **(2) Transparency-Performance Trade-off:** Highly capable black-box models often sacrifice interpretability, whereas intrinsically transparent models may lag in performance. This tension complicates scientific adoption and raises uncertainty about whether novel outputs represent genuine discoveries or the recombination of existing data [475].

10.5. AI for Dynamic and Real-Time Optimized Scientific Experimentation

Real-time AI models can automatically adjust experimental protocols in response to unforeseen variables or shifting conditions, while performing immediate data analysis to substantially enhance research efficiency and innovative potential. Two prominent research directions have emerged: **(1) Agentic Real-Time AI:** This approach advances AI beyond passive data analysis, transforming it into an autonomous research optimizing agent endowed with reasoning, planning, and decision-making capabilities based on real-time experimental feedback. Such agents can systematically survey the literature, generate hypotheses, design experiments, and iteratively refine workflows based on experimental feedback [458, 167, 154]. **(2) Coordination in self-driving laboratories:** These systems integrate robotic platforms, analytical instruments, and AI models into closed-loop frameworks that manage every stage, from experimental planning and execution to data processing. They support applications such as compound screening and novel materials discovery based on real-time signals with minimal human intervention [756, 87, 503].

Despite these advances, two core challenges must be addressed before dynamic, real-time AI experiments become routine: **(1) Reliable integration of heterogeneous devices and AI systems:** Laboratory environments comprise diverse instruments and robotic platforms requiring precise, real-time control and feedback. Systems must ensure compatibility, robustness, and low latency to avoid deviations or downtime caused by integration failures or timing mismatches. **(2) Low-latency decision-making and dynamic optimization:** AI-driven experiments must continuously ingest multisensor and instrument data on the millisecond to second timescales, update model parameters in real-time, and adjust protocols dynamically to maintain workflow continuity and efficiency. Simultaneously, they must uphold robustness and safety to prevent interruptions or hazards due to network jitter or computational bottlenecks [302, 303].

10.6. Multimodal Integration in AI4Research

As scientific data become more diverse, encompassing text, figures, tables, code snippets, and experimental signals, effective multimodal integration has emerged as a linchpin for AI-driven discovery [108, 785, 791, 644, 133]. Early work [245, 121, 534, 675] show that jointly embedding text and figures can substantially boost deep analysis and literature-based discovery, yet this approach often falters when aligning highly specialized diagrams with their textual descriptions [132, 615]. There are two main integration strategies: **(1) Rigorous Multi-Source Data Ingestion:** Scientific datasets span manuscripts, high-resolution images, time-series signals, code artifacts, and structured tables. Each modality requires tailored preprocessing, such as OCR for figures, noise filtering for sensor data, syntax checking for code, to preserve integrity and alignment with domain ontologies [241, 423]. **(2) Interactive Human-in-the-Loop Refinement:** Unlike general-purpose systems, research workflows integrate expert feedback at multiple stages. Interactive interfaces enable domain scientists to validate figure captions, correct table alignments, or adjust the experimental setting based on the multi-modal signals, creating an iterative loop that refines model outputs and builds trust [688, 803, 917].

Nonetheless, multimodal integration in AI4Research faces two core challenges: **(1) Scarcity of cross-modal data and annotation bottleneck:** High-quality aligned annotations are exceedingly scarce, particularly in specialized scientific domains where expert involvement is required for fine-grained pairing, leading to a dramatic escalation of training and evaluation costs. **(2) Quantification of inter-modal uncertainty:** Data originating from diverse sources contain heterogeneous noise; how to uniformly quantify and propagate this uncertainty to support reliable scientific decision-making remains an open challenge.

10.7. Multilingual Integration in AI4Research

Scientific research transcends linguistic and geographic borders. Global initiatives, such as COVID-19 containment and climate modeling, depend on integrating literature, datasets, and expert insights across diverse languages efficiently. If AI tools favor only English or other high-resource languages, research sharing suffers, reinforcing “information silos” and the “knowledge divide” [20, 19]. Most researchers’ native languages lie in the “long tail” of AI systems. Neglecting low-resource languages limits discoverability and citation of high-quality studies and sidelines region-specific topics (e.g., tropical agriculture, minority health). Multilingual pre-training and data augmentation can generate accurate summaries, retrievals, and translations in low-resource languages, breaking down academic barriers [145, 273, 138, 617, 601]. There are two principal integration strategies: **(1) Alignment of Scientific Terminology:** Reproducibility demands consistent terms and semantic fidelity. Multilingual terminology alignment and contextual-fidelity techniques ensure accurate translation of experiments and publications, so researchers worldwide build on a common knowledge base [657, 924, 213]. **(2) Equilibrating Multilingual Performance:** Data imbalances between high- and low-resource languages hinder cross-lingual transfer. Equalizing performance across languages

enhances zero-shot and few-shot capabilities in research applications [138, 617, 792, 907].

Nonetheless, multilingual integration in AI4Research faces two core challenges: **(1) Balancing Capacity and Coverage:** Under finite computational and parameter budgets, striking the right balance between supporting core research capabilities and maintaining broad multilingual performance is critical to prevent “language breadth” from sacrificing “research depth”. This requires fine-grained architectural pruning and resource allocation tailored to specific domains and language pairs. **(2) Analysis of Cross-Lingual Academic Rhetorical Fidelity:** Ensuring that conceptual meanings remain consistent across different languages, preserving the logical integrity of academic argumentation in translation, and addressing language-specific academic conventions constitute important directions for future research.

11. Related work

Recent years have seen increasing interest in AI-assisted or autonomous research across multiple research communities [43]. The empirical use of large language models (LLMs) in research workflows indicates that most researchers are incorporating these models into their processes [464]. Additionally, Yu and Jin [872] survey and predict the rise in AI4Science publications, suggesting strategies to empower AI researchers. Early survey [10, 637, 910, 759, 905] summarize how LLMs are transforming scientific discovery [184, 929, 252, 941]. Li et al. [443] focus more on the ideation developments for LLM-assisted ideation, while Kulkarni et al. [392] and Ren et al. [639] summarize the architectures and benchmarks for LLM-driven discovery methods. Chen et al. [112] propose the Science-of-Science framework, which surveys the AI4Science in multi-agent simulation perspective. Meanwhile, Huang et al. [318] describe the AI-driven scientific discovery process from the perspective of the hypothesis lifecycle [332]. In particular, Zhou et al. [935] and Luo et al. [517] develop a three-stage taxonomy to systematically review assistance role in each phase. Building on this framework, Alkan et al. [15] and Bazgir et al. [51] offer a comprehensive classification of LLM-based hypothesis generation methods. In response to the peer-review crisis, Kim et al. [374] focus more on the bidirectional feedback system with certified reviewers, while Bolanos et al. [69] and Zhuang et al. [946] review the rise of automated scientific paper reviews, which coexist with human oversight.

Although significant advancements have been made in AI4Research, much of the existing survey has focused primarily on scientific discovery and academic writing, often under the umbrella of AI4Science or the limited research stages. However, these discussions typically overlook the broader research lifecycle, including scientific comprehension, academic survey, and peer review. Additionally, they tend to neglect AI applications across these stages. This paper introduces the AI4Research framework and offers a systematic survey of key factors and recent developments driving AI-enabled research. Our goal is to provide the research community with streamlined access to essential resources and insights, thereby facilitating innovative breakthroughs.

12. Conclusion

In conclusion, rapid advancements in artificial intelligence, particularly large language models like OpenAI-o1 and DeepSeek-R1, have demonstrated substantial potential in areas such as logical reasoning and experimental coding. These developments have sparked increasing interest in applying AI to scientific research. However, despite the growing potential of AI in this domain, there is a lack of comprehensive surveys that consolidate current knowledge, hindering further progress. This paper addresses this gap by providing a detailed survey and unified framework for AI4Research. Our contributions include a systematic taxonomy for classifying AI4Research tasks, identification of key research gaps and future directions, and a compilation of open-source resources to support the community. We believe this work will enhance our understanding of AI’s role in research and serve as a catalyst for future advancements in the field.

References

- [1] Marah Abdin, Jyoti Aneja, Harkirat Behl, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar, Michael Harrison, Russell J Hewett, Mojan Javaheripi, Piero Kauffmann, et al. Phi-4 technical report. *arXiv preprint arXiv:2412.08905*, 2024.
- [2] Josh Abramson, Jonas Adler, Jack Dunger, Richard Evans, Tim Green, Alexander Pritzel, Olaf Ronneberger, Lindsay Willmore, Andrew J Ballard, Joshua Bambrick, et al. Accurate structure prediction of biomolecular interactions with alphafold 3. *Nature*, 630(8016):493–500, May 2024.
- [3] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [4] Pierre Achkar, Tim Gollub, and Martin Potthast. Ask, retrieve, summarize: A modular pipeline for scientific literature summarization. *arXiv preprint arXiv:2505.16349*, 2025.
- [5] Steven Adriaensen, Herilalaina Rakotoarison, Samuel Müller, and Frank Hutter. Efficient bayesian learning curve extrapolation using prior-data fitted networks. *Advances in Neural Information Processing Systems*, 36:19858–19886, Dec 2023.
- [6] Tejumade Afonja, Ivaxi Sheth, Ruta Binkyte, Waqar Hanif, Thomas Ulas, Matthias Becker, and Mario Fritz. Llm4grn: Discovering causal gene regulatory networks with llms—evaluation through synthetic data generation. *arXiv preprint arXiv:2410.15828*, 2024.
- [7] Shubham Agarwal, Gaurav Sahu, Abhay Puri, Issam H Laradji, Krishnamurthy DJ Dvijotham, Jason Stanley, Laurent Charlin, and Christopher Pal. Litllm: A toolkit for scientific literature review. *arXiv preprint arXiv:2402.01788*, 2024.
- [8] Shubham Agarwal, Gaurav Sahu, Abhay Puri, Issam H Laradji, Krishnamurthy DJ Dvijotham, Jason Stanley, Laurent Charlin, and Christopher Pal. Llms for literature review: Are we there yet? *arXiv preprint arXiv:2412.15249*, 2024.
- [9] Shubham Agarwal, Gaurav Sahu, Abhay Puri, Issam H. Laradji, Krishnamurthy Dj Dvijotham, Jason Stanley, Laurent Charlin, and Christopher Pal. LitLLMs, LLMs for literature review: Are we there yet? *Transactions on Machine Learning Research*, Apr 2025. ISSN 2835-8856. URL <https://openreview.net/forum?id=heeJqQXKg7>.
- [10] Ajay Agrawal, John McHale, and Alexander Oettl. Artificial intelligence and scientific discovery: A model of prioritized search. *Research Policy*, 53(5):104989, Jun 2024.
- [11] Melih Agraz, Dincer Goksuluk, Peng Zhang, Bum-Rak Choi, Richard T Clements, Gaurav Choudhary, and George Em Karniadakis. Ml-gap: machine learning-enhanced genomic analysis pipeline using autoencoders and data augmentation. *Frontiers in Genetics*, 15:1442759, Sep 2024.
- [12] Intology AI. Zochi technical report, Mar 2025. URL https://github.com/IntologyAI/Zochi/blob/main/Zochi_Technical_Report.pdf. Zochi Technical Report.
- [13] Nurmukhammed Aitymbetov and Dimitrios Zorbas. Autonomous machine learning-based peer reviewer selection system. In *Proceedings of the 31st International Conference on Computational Linguistics: System Demonstrations*, pages 199–207, Jan 2025.

- [14] Andres Algaba, Vincent Holst, Floriano Tori, Melika Mobini, Brecht Verbeken, Sylvia Wenmackers, and Vincent Ginis. How deep do large language models internalize scientific literature and citation practices? *arXiv preprint arXiv:2504.02767*, 2025.
- [15] Atilla Kaan Alkan, Shashwat Sourav, Maja Jablonska, Simone Astarita, Rishabh Chakrabarty, Nikhil Garuda, Pranav Khetarpal, Maciej Pióro, Dimitrios Tanoglidis, Kartheik G Iyer, et al. A survey on hypothesis generation for scientific discovery in the era of large language models. *arXiv preprint arXiv:2504.05496*, 2025.
- [16] Enes Altuncu, Jason RC Nurse, Meryem Bagriacik, Sophie Kaleba, Haiyue Yuan, Lisa Bonheme, and Shujun Li. aedfact: Scientific fact-checking made easier via semi-automatic discovery of relevant expert opinions. *arXiv preprint arXiv:2305.07796*, 2023.
- [17] Carlos Alvarez, Maxwell Bennett, and Lucy Lu Wang. Zero-shot scientific claim verification using llms and citation text. In *Proceedings of the Fourth Workshop on Scholarly Document Processing (SDP 2024)*, pages 269–276, Aug 2024.
- [18] Jose M Alvarez, Alejandra Bringas Colmenarejo, Alaa Elobaid, Simone Fabbriizzi, Miriam Fahimi, Antonio Ferrara, Siamak Ghodsi, Carlos Mougán, Ioanna Papageorgiou, Paula Reyero, et al. Policy advice and best practices on bias and fairness in ai. *Ethics and Information Technology*, 26(2):31, Apr 2024.
- [19] Tatsuya Amano, Juan P González-Varo, and William J Sutherland. Languages are still a major barrier to global science. *PLoS biology*, 14(12):e2000933, Dec 2016.
- [20] Tatsuya Amano, Clarissa Rios Rojas, Yap Boum II, Margarita Calvo, and Biswapriya B Misra. Ten tips for overcoming language barriers in science. *Nature Human Behaviour*, 5(9):1119–1122, Jul 2021.
- [21] Barrett R Anderson, Jash Hemant Shah, and Max Kreminski. Homogenization effects of large language models on human creative ideation. In *Proceedings of the 16th conference on creativity & cognition*, pages 413–425, Jun 2024.
- [22] Nicholas H Angello, David M Friday, Changhyun Hwang, Seungjoo Yi, Austin H Cheng, Tiara C Torres-Flores, Edward R Jira, Wesley Wang, Alán Aspuru-Guzik, Martin D Burke, et al. Closed-loop transfer enables artificial intelligence to yield chemical knowledge. *Nature*, 633(8029):351–358, 2024.
- [23] Angelos Angelopoulos, James F Cahoon, and Ron Alterovitz. Transforming science labs into automated factories of discovery. *Science Robotics*, 9(95):eadm6991, 2024.
- [24] AI Anthropic. The claude 3 model family: Opus, sonnet, haiku. Claude-3 Model Card, Mar 2024.
- [25] Sören Arlt, Haonan Duan, Felix Li, Sang Michael Xie, Yuhuai Wu, and Mario Krenn. Meta-designing quantum experiments with language models. *arXiv preprint arXiv:2406.02470*, 2024.
- [26] Akari Asai, Jacqueline He, Rulin Shao, Weijia Shi, Amanpreet Singh, Joseph Chee Chang, Kyle Lo, Luca Soldaini, Sergey Feldman, Mike D’arcy, et al. Openscholar: Synthesizing scientific literature with retrieval-augmented lms. *arXiv preprint arXiv:2411.14199*, 2024.
- [27] Joshua Ashkinaze, Julia Mendelsohn, Li Qiwei, Ceren Budak, and Eric Gilbert. How ai ideas affect the creativity, diversity, and evolution of human ideas: evidence from a large, dynamic experiment. *arXiv preprint arXiv:2401.13481*, 2024.

- [28] Shir Ashury-Tahan, Yifan Mai, Ariel Gera, Yotam Perlitz, Asaf Yehudai, Elron Bandel, Leshem Choshen, Eyal Shnarch, Percy Liang, Michal Shmueli-Scheuer, et al. The mighty torr: A benchmark for table reasoning and robustness. *arXiv preprint arXiv:2502.19412*, 2025.
- [29] Assafelovic. gpt-researcher, May 2023. URL <https://github.com/assafelovic/gpt-researcher>. gpt-researcher.
- [30] Pepa Atanasova. Generating fact checking explanations. In *Accountable and Explainable Methods for Complex Reasoning over Text*, pages 83–103. Springer, Apr 2024.
- [31] Steven Au, Cameron J Dimacali, Ojasmitha Pedirappagari, Namyong Park, Franck Dernoncourt, Yu Wang, Nikos Kanakaris, Hanieh Deilamsalehy, Ryan A Rossi, and Nesreen K Ahmed. Personalized graph-based retrieval for large language models. *arXiv preprint arXiv:2501.02157*, 2025.
- [32] Sören Auer, Dante AC Barone, Cassiano Bartz, Eduardo G Cortes, Mohamad Yaser Jaradeh, Oliver Karras, Manolis Koubarakis, Dmitry Mouromtsev, Dmitrii Pliukhin, Daniil Radyush, et al. The sciqa scientific question answering benchmark for scholarly knowledge. *Scientific Reports*, 13(1):7240, May 2023.
- [33] Axios. Self-driving labs are the new ai asset. *Axios*, Aug 2024. URL <https://www.axios.com/2024/08/09/ai-self-driving-science-labs-research>.
- [34] Shingo Ayabe, Takuto Otomo, Hiroshi Kera, and Kazuhiko Kawamoto. Robustness evaluation of offline reinforcement learning for robot control against action perturbations. *arXiv preprint arXiv:2412.18781*, 2024.
- [35] Ibrahim Al Azher, Miftahul Jannat Mokarrama, Zhishuai Guo, Sagnik Ray Choudhury, and Hamed Alhoori. Futuregen: Llm-rag approach to generate the future work of scientific article. *arXiv preprint arXiv:2503.16561*, 2025.
- [36] 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.
- [37] Xiaomei Bai, Mengyang Wang, Ivan Lee, Zhuo Yang, Xiangjie Kong, and Feng Xia. Scientific paper recommendation: A survey. *Ieee Access*, 7:9324–9339, Jan 2019.
- [38] Honglin Bao, Siyang Wu, Jiwoong Choi, Yingrong Mao, and James A Evans. Language models surface the unwritten code of science and society. *arXiv preprint arXiv:2505.18942*, 2025.
- [39] Zhijie Bao, Qingyun Liu, Ying Guo, Zhengqiang Ye, Jun Shen, Shirong Xie, Jiajie Peng, Xuanjing Huang, and Zhongyu Wei. Piors: Personalized intelligent outpatient reception based on large language model with multi-agents medical scenario simulation. *arXiv preprint arXiv:2411.13902*, 2024.
- [40] Mitra Baratchi, Can Wang, Steffen Limmer, Jan N van Rijn, Holger Hoos, Thomas Bäck, and Markus Olhofer. Automated machine learning: past, present and future. *Artificial intelligence review*, 57(5): 122, Apr 2024.
- [41] Gregory Barber. Google deepmind’s ai dreamed up 380,000 new materials. the next challenge is making them. *WIRED*, Nov 2023. URL <https://www.wired.com/story/an-ai-dreamed-up-380000-new-materials-the-next-challenge-is-making-them/>.

- [42] Rafael Barbudo, Sebastián Ventura, and José Raúl Romero. Eight years of automl: categorisation, review and trends. *Knowledge and Information Systems*, 65(12):5097–5149, Aug 2023.
- [43] Kristian G Barman, Sascha Caron, Emily Sullivan, Henk W de Regt, Roberto Ruiz de Austri, Mieke Boon, Michael Färber, Stefan Fröse, Faegheh Hasibi, Andreas Ipp, et al. Large physics models: Towards a collaborative approach with large language models and foundation models. *arXiv preprint arXiv:2501.05382*, 2025.
- [44] Annabel R Basford, Aaron H Bernardino, Paula CP Teeuwen, Benjamin D Egleston, Joshua Humphreys, Kim E Jelfs, Jonathan R Nitschke, Imogen A Riddell, and Rebecca L Greenaway. Development of an automated workflow for screening the assembly and host–guest behavior of metal-organic cages towards accelerated discovery. *Angewandte Chemie International Edition*, page e202424270, Apr 2024.
- [45] Setio Basuki and Masatoshi Tsuchiya. The quality assist: A technology-assisted peer review based on citation functions to predict the paper quality. *IEEE Access*, 10:126815–126831, Dec 2022.
- [46] Bruno C Batista, SV Amrutha, Jie Yan, Beni B Dangi, and Oliver Steinbock. High-throughput robotic collection, imaging, and machine learning analysis of salt patterns: composition and concentration from dried droplet photos. *Digital Discovery*, 4(4):1030–1041, Feb 2025.
- [47] Peter Battaglia, Razvan Pascanu, Matthew Lai, Danilo Jimenez Rezende, et al. Interaction networks for learning about objects, relations and physics. *Advances in Neural Information Processing Systems*, 29, Dec 2016.
- [48] Tim Baumgärtner, Ted Briscoe, and Iryna Gurevych. Peerqa: A scientific question answering dataset from peer reviews. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 508–544, Feb 2025.
- [49] Atılım Güneş Baydin, Kyle Cranmer, Pablo de Castro Manzano, Christophe Delaere, Denis Derkach, Julien Donini, Tommaso Dorigo, Andrea Giammanco, Jan Kieseler, Lukas Layer, et al. Toward machine learning optimization of experimental design. *Nuclear Physics News*, 31(1):25–28, Feb 2021.
- [50] Adrián Bazaga, Pietro Lio, and Gos Micklem. Unsupervised pretraining for fact verification by language model distillation. *arXiv preprint arXiv:2309.16540*, 2023.
- [51] Adib Bazgir, Yuwen Zhang, et al. Agentichypothesis: A survey on hypothesis generation using llm systems. *Towards Agentic AI for Science: Hypothesis Generation, Comprehension, Quantification, and Validation*, Mar 2025.
- [52] Joeran Beel, Bela Gipp, Stefan Langer, and Corinna Breitingner. Paper recommender systems: a literature survey. *International Journal on Digital Libraries*, 17(4):305–338, Jul 2016.
- [53] Joeran Beel, Min-Yen Kan, and Moritz Baumgart. Evaluating sakana’s ai scientist for autonomous research: Wishful thinking or an emerging reality towards’ artificial research intelligence’(ari)? *arXiv preprint arXiv:2502.14297*, 2025.
- [54] Jonas Belouadi, Anne Lauscher, and Steffen Eger. Automatizkz: Text-guided synthesis of scientific vector graphics with tikz. *arXiv preprint arXiv:2310.00367*, 2023.

- [55] Jonas Belouadi, Eddy Ilg, Margret Keuper, Hideki Tanaka, Masao Utiyama, Raj Dabre, Steffen Eger, and Simone Paolo Ponzetto. Tikzero: Zero-shot text-guided graphics program synthesis. *arXiv preprint arXiv:2503.11509*, 2025.
- [56] Iz Beltagy, Kyle Lo, and Arman Cohan. SciBERT: A pretrained language model for scientific text. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan, editors, *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, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1371. URL <https://aclanthology.org/D19-1371/>.
- [57] Leonard Bereska and Efstratios Gavves. Mechanistic interpretability for ai safety—a review. *arXiv preprint arXiv:2404.14082*, 2024.
- [58] Marialena Bevilacqua, Kezia Oketch, Ruiyang Qin, Will Stamey, Xinyuan Zhang, Yi Gan, Kai Yang, and Ahmed Abbasi. When automated assessment meets automated content generation: Examining text quality in the era of gpts. *ACM Trans. Inf. Syst.*, 43(2), January 2025. ISSN 1046-8188. doi: 10.1145/3702639. URL <https://doi.org/10.1145/3702639>.
- [59] Prabhat Kumar Bharti, Shashi Ranjan, Tirthankar Ghosal, Mayank Agrawal, and Asif Ekbal. Peerassist: leveraging on paper-review interactions to predict peer review decisions. In *Towards Open and Trustworthy Digital Societies: 23rd International Conference on Asia-Pacific Digital Libraries, ICADL 2021, Virtual Event, December 1–3, 2021, Proceedings 23*, pages 421–435. Springer, Nov 2021.
- [60] Prabhat Kumar Bharti, Meith Navlakha, Mayank Agarwal, and Asif Ekbal. Politepeer: does peer review hurt? a dataset to gauge politeness intensity in the peer reviews. *Language Resources and Evaluation*, 58(4):1291–1313, May 2024.
- [61] Haiyang Bian, Yixin Chen, Erpai Luo, Xinze Wu, Minsheng Hao, Lei Wei, and Xuegong Zhang. General-purpose pre-trained large cellular models for single-cell transcriptomics. *National Science Review*, 11(11):nwae340, Sep 2024.
- [62] Junyi Bian, Xiaolei Qin, Wuhe Zou, Mengzuo Huang, Congyi Luo, Ke Zhang, and Weidong Zhang. Helm: Highlighted evidence augmented language model for enhanced table-to-text generation. *arXiv preprint arXiv:2311.08896*, 2023.
- [63] Thulasi Bikku, Nirmala Rani Narimalla, Keerthi Konda, Anusha Nakkala, Avanti Yarlagadda, and B Sachuthanathan. Generating accurate and engaging research paper titles using nlp techniques. In *International Conference on Innovations in Bio-Inspired Computing and Applications*, pages 428–437. Springer, May 2025.
- [64] Marcel Binz and Eric Schulz. Using cognitive psychology to understand gpt-3. *Proceedings of the National Academy of Sciences*, 120(6), February 2023. ISSN 1091-6490. doi: 10.1073/pnas.2218523120. URL <http://dx.doi.org/10.1073/pnas.2218523120>.
- [65] Marah Blaurock, Marion Büttgen, and Jeroen Schepers. Designing collaborative intelligence systems for employee-ai service co-production. *Journal of Service Research*, page 10946705241238751, Mar 2024.
- [66] Severin Bochem, Eduardo Gonzalez-Sanchez, Yves Bicker, and Gabriele Fadini. Improving generalization of robot locomotion policies via sharpness-aware reinforcement learning. *arXiv preprint arXiv:2411.19732*, 2024.

- [67] Amber Boehnlein, Markus Diefenthaler, Nobuo Sato, Malachi Schram, Veronique Ziegler, Cristiano Fanelli, Morten Hjorth-Jensen, Tanja Horn, Michelle P Kuchera, Dean Lee, et al. Colloquium: Machine learning in nuclear physics. *Reviews of modern physics*, 94(3):031003, Sep 2022.
- [68] Daniil A Boiko, Robert MacKnight, Ben Kline, and Gabe Gomes. Autonomous chemical research with large language models. *Nature*, 624(7992):570–578, Dec 2023.
- [69] Francisco Bolanos, Angelo Salatino, Francesco Osborne, and Enrico Motta. Artificial intelligence for literature reviews: Opportunities and challenges. *Artificial Intelligence Review*, 57(10):259, Aug 2024.
- [70] Valeriia Bolotova, Vladislav Blinov, Falk Scholer, W Bruce Croft, and Mark Sanderson. A non-factoid question-answering taxonomy. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 1196–1207, Jul 2022.
- [71] Elliot Bolton, Abhinav Venigalla, Michihiro Yasunaga, David Hall, Betty Xiong, Tony Lee, Roxana Daneshjou, Jonathan Frankle, Percy Liang, Michael Carbin, et al. Biomedlm: A 2.7 b parameter language model trained on biomedical text. *arXiv preprint arXiv:2403.18421*, 2024.
- [72] Necva Bölücü, Yunus Can Bilge, Dilber Çetintaş, and Zehra Yücel. Modest: A dataset for multi domain scientific title generation. *Knowledge-Based Systems*, page 113557, Jun 2025.
- [73] Jan A Botha, Manaal Faruqui, John Alex, Jason Baldridge, and Dipanjan Das. Learning to split and rephrase from wikipedia edit history. *arXiv preprint arXiv:1808.09468*, 2018.
- [74] Jack Boylan, Shashank Mangla, Dominic Thorn, Demian Gholipour Ghalandari, Parsa Ghaffari, and Chris Hokamp. Kgvalidator: A framework for automatic validation of knowledge graph construction. *arXiv preprint arXiv:2404.15923*, 2024.
- [75] Tobias Braun, Mark Rothermel, Marcus Rohrbach, and Anna Rohrbach. Defame: Dynamic evidence-based fact-checking with multimodal experts. *arXiv preprint arXiv:2412.10510*, 2024.
- [76] Thomas Britton, Cullan Bedwell, Abhijeet Chawhan, Julie Crowe, Naomi Jarvis, Torri Jeske, Nikhil Kalra, David Lawrence, and Diana McSpadden. Ai driven experiment calibration and control. In *EPJ Web of Conferences*, volume 295, page 02003. EDP Sciences, May 2024.
- [77] Victor Brodsky, Ehsan Ullah, Andrey Bychkov, Andrew H Song, Eric E Walk, Peter Louis, Ghulam Rasool, Rajendra S Singh, Faisal Mahmood, Marilyn M Bui, et al. Generative artificial intelligence in anatomic pathology. *Archives of Pathology & Laboratory Medicine*, Apr 2025.
- [78] Qingwen Bu, Jia Zeng, Li Chen, Yanchao Yang, Guyue Zhou, Junchi Yan, Ping Luo, Heming Cui, Yi Ma, and Hongyang Li. Closed-loop visuomotor control with generative expectation for robotic manipulation. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [79] Markus J Buehler. Accelerating scientific discovery with generative knowledge extraction, graph-based representation, and multimodal intelligent graph reasoning. *Machine Learning: Science and Technology*, 5(3):035083, Sep 2024.
- [80] Lukas Buess, Matthias Keicher, Nassir Navab, Andreas Maier, and Soroosh Tayebi Arasteh. From large language models to multimodal ai: A scoping review on the potential of generative ai in medicine. *arXiv preprint arXiv:2502.09242*, 2025.

- [81] Charlotte Bunne, Yusuf Roohani, Yanay Rosen, Ankit Gupta, Xikun Zhang, Marcel Roed, Theo Alexandrov, Mohammed AlQuraishi, Patricia Brennan, Daniel B Burkhardt, et al. How to build the virtual cell with artificial intelligence: Priorities and opportunities. *Cell*, 187(25):7045–7063, Dec 2024.
- [82] James Burgess, Jeffrey J Nirschl, Laura Bravo-Sánchez, Alejandro Lozano, Sanket Rajan Gupte, Jesus G Galaz-Montoya, Yuhui Zhang, Yuchang Su, Disha Bhowmik, Zachary Coman, et al. Microvqa: A multimodal reasoning benchmark for microscopy-based scientific research. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 19552–19564, Mar 2025.
- [83] Keith T Butler, Daniel W Davies, Hugh Cartwright, Olexandr Isayev, and Aron Walsh. Machine learning for molecular and materials science. *Nature*, 559(7715):547–555, Jul 2018.
- [84] Necva Bölücü, Yunus Can Bilge, Dilber Çetintaş, and Zehra Yücel. Modest: A dataset for multi domain scientific title generation. *Knowledge-Based Systems*, 321:113557, Jun 2025. ISSN 0950-7051. doi: <https://doi.org/10.1016/j.knosys.2025.113557>. URL <https://www.sciencedirect.com/science/article/pii/S0950705125006033>.
- [85] Yi Cai and Gerhard Wunder. On gradient-like explanation under a black-box setting: when black-box explanations become as good as white-box. *arXiv preprint arXiv:2308.09381*, 2023.
- [86] Reyes Calderon and Francisco Herrera. And plato met chatgpt: an ethical reflection on the use of chatbots in scientific research writing, with a particular focus on the social sciences. *Humanities and Social Sciences Communications*, 12(1):1–13, May 2025.
- [87] Richard B Canty, Jeffrey A Bennett, Keith A Brown, Tonio Buonassisi, Sergei V Kalinin, John R Kitchin, Benji Maruyama, Robert G Moore, Joshua Schrier, Martin Seifrid, et al. Science acceleration and accessibility with self-driving labs. *Nature Communications*, 16(1):3856, Apr 2025.
- [88] Lang Cao and Hanbing Liu. Tablemaster: A recipe to advance table understanding with language models. *arXiv preprint arXiv:2501.19378*, 2025.
- [89] Shuxiang Cao, Zijian Zhang, Mohammed Alghadeer, Simone D Fasciati, Michele Piscitelli, Mustafa Bakr, Peter Leek, and Alán Aspuru-Guzik. Agents for self-driving laboratories applied to quantum computing. *arXiv preprint arXiv:2412.07978*, 2024.
- [90] Stanley Cao and Kevin Liu. Figuring out figures: Using textual references to caption scientific figures. *arXiv preprint arXiv:2407.11008*, 2024.
- [91] Yupeng Cao, Aishwarya Muralidharan Nair, Elyon Eyimife, Nastaran Jamalipour Soofi, KP Subbalakshmi, John R Wullert II, Chumki Basu, and David Shallcross. Can large language models detect misinformation in scientific news reporting? *arXiv preprint arXiv:2402.14268*, 2024.
- [92] Ege Yiğit Çelik and Selma Tekir. Citebart: Learning to generate citations for local citation recommendation. *arXiv preprint arXiv:2412.17534*, 2024.
- [93] Tuhin Chakrabarty, Philippe Laban, Divyansh Agarwal, Smaranda Muresan, and Chien-Sheng Wu. Art or artifice? large language models and the false promise of creativity. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, pages 1–34, May 2024.

- [94] Eric Chamoun, Michael Schlichtkrull, and Andreas Vlachos. Automated focused feedback generation for scientific writing assistance. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Findings of the Association for Computational Linguistics: ACL 2024*, pages 9742–9763, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.580. URL <https://aclanthology.org/2024.findings-acl.580/>.
- [95] Jun Shern Chan, Neil Chowdhury, Oliver Jaffe, James Aung, Dane Sherburn, Evan Mays, Giulio Starace, Kevin Liu, Leon Maksin, Tejal Patwardhan, et al. Mle-bench: Evaluating machine learning agents on machine learning engineering. *arXiv preprint arXiv:2410.07095*, 2024.
- [96] Mohit Chandra, Suchismita Naik, Denae Ford, Ebele Okoli, Munmun De Choudhury, Mahsa Ershadi, Gonzalo Ramos, Javier Hernandez, Ananya Bhattacharjee, Shahed Warreth, et al. From lived experience to insight: Unpacking the psychological risks of using ai conversational agents. *arXiv preprint arXiv:2412.07951*, 2024.
- [97] Yuan Chang, Ziyue Li, Hengyuan Zhang, Yuanbo Kong, Yanru Wu, Zhijiang Guo, and Ngai Wong. Treereview: A dynamic tree of questions framework for deep and efficient llm-based scientific peer review. *arXiv preprint arXiv:2506.07642*, 2025.
- [98] Laurent Charlin and Richard Zemel. The toronto paper matching system: an automated paper-reviewer assignment system. May 2013.
- [99] Laurent Charlin, Richard Zemel, and Craig Boutilier. A framework for optimizing paper matching. In *Proceedings of the Twenty-Seventh Conference on Uncertainty in Artificial Intelligence*, pages 86–95, Jul 2011.
- [100] Rajdip Chaudhuri. Llm based exploratory data analysis using bigquery data canvas, Oct 2024. URL <https://medium.com/google-cloud/llm-based-exploratory-data-analysis-using-bigquery-data-canvas-42fbecb9f009>. LLM Based Exploratory Data Analysis Using BigQuery Data Canvas.
- [101] Chi Chen, Weike Ye, Yunxing Zuo, Chen Zheng, and Shyue Ping Ong. Graph networks as a universal machine learning framework for molecules and crystals. *Chemistry of Materials*, 31(9):3564–3572, Apr 2019.
- [102] Hui Chen, Miao Xiong, Yujie Lu, Wei Han, Ailin Deng, Yufei He, Jiaying Wu, Yibo Li, Yue Liu, and Bryan Hooi. Mlr-bench: Evaluating ai agents on open-ended machine learning research. *arXiv preprint arXiv:2505.19955*, 2025.
- [103] Jingqiang Chen and Hai Zhuge. Automatic generation of related work through summarizing citations. *Concurrency and Computation: Practice and Experience*, 31(3):e4261, Sep 2019.
- [104] Junlan Chen, Kexin Zhang, Daifeng Li, Yangyang Feng, Yuxuan Zhang, and Bowen Deng. Structuring scientific innovation: A framework for modeling and discovering impactful knowledge combinations. *arXiv preprint arXiv:2503.18865*, 2025.
- [105] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.

- [106] Nuo Chen, Andre Lin HuiKai, Jiaying Wu, Junyi Hou, Zining Zhang, Qian Wang, Xidong Wang, and Bingsheng He. Xtragpt: Llms for human-ai collaboration on controllable academic paper revision. *arXiv preprint arXiv:2505.11336*, 2025.
- [107] Qiguang Chen, Libo Qin, Jiaqi Wang, Jingxuan Zhou, and Wanxiang Che. Unlocking the capabilities of thought: A reasoning boundary framework to quantify and optimize chain-of-thought. *Advances in Neural Information Processing Systems*, 37:54872–54904, Sep 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/hash/62ab1c2cb4b03e717005479efb211841-Abstract-Conference.html.
- [108] Qiguang Chen, Libo Qin, Jin Zhang, Zhi Chen, Xiao Xu, and Wanxiang Che. M³CoT: A novel benchmark for multi-domain multi-step multi-modal chain-of-thought. pages 8199–8221, August 2024. doi: 10.18653/v1/2024.acl-long.446. URL <https://aclanthology.org/2024.acl-long.446/>.
- [109] Qiguang Chen, Libo Qin, Jinhao Liu, Yue Liao, Jiaqi Wang, Jingxuan Zhou, and Wanxiang Che. Rbf++: Quantifying and optimizing reasoning boundaries across measurable and unmeasurable capabilities for chain-of-thought reasoning. *arXiv preprint arXiv:2505.13307*, 2025.
- [110] Qiguang Chen, Libo Qin, Jinhao Liu, Dengyun Peng, Jiannan Guan, Peng Wang, Mengkang Hu, Yuhang Zhou, Te Gao, and Wangxiang Che. Towards reasoning era: A survey of long chain-of-thought for reasoning large language models. *arXiv preprint arXiv:2503.09567*, 2025.
- [111] Qiguang Chen, Libo Qin, Jinhao Liu, Dengyun Peng, Jiaqi Wang, Mengkang Hu, Zhi Chen, Wanxiang Che, and Ting Liu. Ecm: A unified electronic circuit model for explaining the emergence of in-context learning and chain-of-thought in large language model. *arXiv preprint arXiv:2502.03325*, 2025.
- [112] Renqi Chen, Haoyang Su, Shixiang Tang, Zhenfei Yin, Qi Wu, Hui Li, Ye Sun, Nanqing Dong, Wanli Ouyang, and Philip Torr. Ai-driven automation can become the foundation of next-era science of science research. *arXiv preprint arXiv:2505.12039*, 2025.
- [113] Wei Chen, Han Ding, Meng Yuan, Zhao Zhang, Deqing Wang, and Fuzhen Zhuang. Bridging social psychology and llm reasoning: Conflict-aware meta-review generation via cognitive alignment. *arXiv preprint arXiv:2503.13879*, 2025.
- [114] Wenhui Chen, Ming Yin, Max Ku, Pan Lu, Yixin Wan, Xueguang Ma, Jianyu Xu, Xinyi Wang, and Tony Xia. Theoremqa: A theorem-driven question answering dataset. *arXiv preprint arXiv:2305.12524*, 2023.
- [115] Xiaodong Chen, Jillian M Buriak, Mathieu Salanne, and Huolin Xin. Nano & ai: A nobel partnership, Nov 2024.
- [116] Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. Teaching large language models to self-debug. *arXiv preprint arXiv:2304.05128*, 2023.
- [117] Xiuying Chen, Hind Alamro, Mingzhe Li, Shen Gao, Xiangliang Zhang, Dongyan Zhao, and Rui Yan. Capturing relations between scientific papers: An abstractive model for related work section generation. *Association for Computational Linguistics*, Aug 2021.
- [118] Xiuying Chen, Hind Alamro, Mingzhe Li, Shen Gao, Rui Yan, Xin Gao, and Xiangliang Zhang. Target-aware abstractive related work generation with contrastive learning. In *Proceedings of the 45th international ACM SIGIR conference on research and development in information retrieval*, pages 373–383, Jul 2022.

- [119] Xiuying Chen, Tairan Wang, Taicheng Guo, Kehan Guo, Juexiao Zhou, Haoyang Li, Mingchen Zhuge, Jürgen Schmidhuber, Xin Gao, and Xiangliang Zhang. Scholarchemqa: Unveiling the power of language models in chemical research question answering. *arXiv preprint arXiv:2407.16931*, 2024.
- [120] Yaoyu Chen, Yuheng Hu, and Yingda Lu. Predicting field experiments with large language models. *arXiv preprint arXiv:2504.01167*, 2025.
- [121] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. Uniter: Universal image-text representation learning. In *European conference on computer vision*, pages 104–120. Springer, Sep 2020.
- [122] Ying-Jung Chen, Ahmad Albarqawi, and Chi-Sheng Chen. Reinforcing clinical decision support through multi-agent systems and ethical ai governance. *arXiv preprint arXiv:2504.03699*, 2025.
- [123] Yiqun Chen and James Zou. Genept: a simple but effective foundation model for genes and cells built from chatgpt. *bioRxiv*, pages 2023–10, Mar 2024.
- [124] Yiting Chen, Tracy Xiao Liu, You Shan, and Songfa Zhong. The emergence of economic rationality of gpt. *Proceedings of the National Academy of Sciences*, 120(51):e2316205120, Dec 2023.
- [125] Zhi Chen, Qiguang Chen, Libo Qin, Qipeng Guo, Haijun Lv, Yicheng Zou, Wanxiang Che, Hang Yan, Kai Chen, and Dahua Lin. What are the essential factors in crafting effective long context multi-hop instruction datasets? insights and best practices. *arXiv preprint arXiv:2409.01893*, 2024.
- [126] Zirui Chen, Shijie Chen, Yuting Ning, Qianheng Zhang, Boshi Wang, Botao Yu, Yifei Li, Zeyi Liao, Chen Wei, Zitong Lu, et al. Scienceagentbench: Toward rigorous assessment of language agents for data-driven scientific discovery. *arXiv preprint arXiv:2410.05080*, 2024.
- [127] Adam Cheng, Aaron Calhoun, and Gabriel Reedy. Artificial intelligence-assisted academic writing: recommendations for ethical use. *Advances in Simulation*, 10(1):22, Apr 2025.
- [128] Lu Cheng, Ahmadreza Mosallanezhad, Paras Sheth, and Huan Liu. Causal learning for socially responsible ai. *arXiv preprint arXiv:2104.12278*, 2021.
- [129] Xiaoqing Cheng, Ruizhe Chen, Hongying Zan, Yuxiang Jia, and Min Peng. Biasfilter: An inference-time debiasing framework for large language models. *arXiv preprint arXiv:2505.23829*, 2025.
- [130] Xusen Cheng and Lulu Zhang. Ai-generated literature reviews threaten scientific progress. *Nature*, 641(8064):852–852, 2025.
- [131] Zhi-Qi Cheng, Qi Dai, and Alexander G Hauptmann. Chartreader: A unified framework for chart derendering and comprehension without heuristic rules. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 22202–22213, Apr 2023.
- [132] Zihui Cheng, Qiguang Chen, Xiao Xu, Jiaqi Wang, Weiyun Wang, Hao Fei, Yidong Wang, Alex Jinpeng Wang, Zhi Chen, Wanxiang Che, et al. Visual thoughts: A unified perspective of understanding multimodal chain-of-thought. *arXiv preprint arXiv:2505.15510*, 2025.
- [133] Zihui Cheng, Qiguang Chen, Jin Zhang, Hao Fei, Xiaocheng Feng, Wanxiang Che, Min Li, and Libo Qin. Comt: A novel benchmark for chain of multi-modal thought on large vision-language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 23678–23686, Apr 2025.

- [134] Jaewoong Choi and Byungju Lee. Accelerating materials language processing with large language models. *Communications Materials*, 5(1):13, Feb 2024.
- [135] Zheng Chu, Huiming Fan, Jingchang Chen, Qianyu Wang, Mingda Yang, Jiafeng Liang, Zhongjie Wang, Hao Li, Guo Tang, Ming Liu, et al. Self-critique guided iterative reasoning for multi-hop question answering. *arXiv preprint arXiv:2505.19112*, 2025.
- [136] Zhumin Chu, Qingyao Ai, Yiteng Tu, Haitao Li, and Yiqun Liu. Automatic large language model evaluation via peer review. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management*, pages 384–393, Oct 2024.
- [137] Zhumin Chu, Qingyao Ai, Yiteng Tu, Haitao Li, and Yiqun Liu. Pre: A peer review based large language model evaluator. *arXiv preprint arXiv:2401.15641*, 2024.
- [138] Lynn Chua, Badih Ghazi, Yangsibo Huang, Pritish Kamath, Ravi Kumar, Pasin Manurangsi, Amer Sinha, Chulin Xie, and Chiyuan Zhang. Crosslingual capabilities and knowledge barriers in multilingual large language models. *arXiv preprint arXiv:2406.16135*, 2024.
- [139] Ilker Cingillioglu, Uri Gal, and Artem Prokhorov. Ai-experiments in education: An ai-driven randomized controlled trial for higher education research. *Education and Information Technologies*, 29(15):19649–19677, 2024.
- [140] Defne Circi, Ghazal Khalighinejad, Anlan Chen, Bhuwan Dhingra, and L Catherine Brinson. How well do large language models understand tables in materials science? *Integrating Materials and Manufacturing Innovation*, 13(3):669–687, Jul 2024. URL <https://link.springer.com/article/10.1007/s40192-024-00362-6>.
- [141] Davide Cirillo, Iker Núñez-Carpintero, and Alfonso Valencia. Artificial intelligence in cancer research: learning at different levels of data granularity. *Molecular oncology*, 15(4):817–829, Apr 2021.
- [142] Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. Boolq: Exploring the surprising difficulty of natural yes/no questions. *arXiv preprint arXiv:1905.10044*, 2019.
- [143] Andy Coenen, Luke Davis, Daphne Ippolito, Emily Reif, and Ann Yuan. Wordcraft: A human-ai collaborative editor for story writing. *arXiv preprint arXiv:2107.07430*, 2021.
- [144] Kai-Hendrik Cohrs, Emiliano Diaz, Vasileios Sitokonstantinou, Gherardo Varando, and Gustau Camps-Valls. Large language models for causal hypothesis generation in science. *Machine Learning: Science and Technology*, 6(1):013001, Jan 2025.
- [145] Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. Unsupervised cross-lingual representation learning at scale. *arXiv preprint arXiv:1911.02116*, 2019.
- [146] Micaela E Consens, Cameron Dufault, Michael Wainberg, Duncan Forster, Mehran Karimzadeh, Hani Goodarzi, Fabian J Theis, Alan Moses, and Bo Wang. Transformers and genome language models. *Nature Machine Intelligence*, pages 1–17, Mar 2025.
- [147] Paulo Henrique Couto, Quang Phuoc Ho, Nageeta Kumari, Benedictus Kent Rachmat, Thanh Gia Hieu Khuong, Ihsan Ullah, and Lisheng Sun-Hosoya. Relevai-reviewer: A benchmark on ai reviewers for survey paper relevance. *arXiv preprint arXiv:2406.10294*, 2024.

- [148] Douglas B Craig. A human-LLM note-taking system with case-based reasoning as framework for scientific discovery. In Peter Jansen, Bhavana Dalvi Mishra, Harsh Trivedi, Bodhisattwa Prasad Majumder, Tom Hope, Tushar Khot, Doug Downey, and Eric Horvitz, editors, *Proceedings of the 1st Workshop on AI and Scientific Discovery: Directions and Opportunities*, pages 22–30, Albuquerque, New Mexico, USA, May 2025. Association for Computational Linguistics. ISBN 979-8-89176-224-4. doi: 10.18653/v1/2025.aisd-main.3. URL <https://aclanthology.org/2025.aisd-main.3/>.
- [149] Miles Cranmer, Sam Greydanus, Stephan Hoyer, Peter Battaglia, David Spergel, and Shirley Ho. Lagrangian neural networks. *arXiv preprint arXiv:2003.04630*, 2020.
- [150] Haotian Cui, Chloe Wang, Hassaan Maan, Kuan Pang, Fengning Luo, Nan Duan, and Bo Wang. scgpt: toward building a foundation model for single-cell multi-omics using generative ai. *Nature Methods*, 21(8):1470–1480, Feb 2024.
- [151] Haotian Cui, Yue Xu, Kuan Pang, Gen Li, Fanglin Gong, Bo Wang, and Bowen Li. Lumi-lab: a foundation model-driven autonomous platform enabling discovery of new ionizable lipid designs for mrna delivery. *BioRxiv*, pages 2025–02, Feb 2025.
- [152] Ziyang Cui, Ning Li, and Huaikang Zhou. Can ai replace human subjects? a large-scale replication of psychological experiments with llms. *arXiv preprint arXiv:2409.00128*, 2024.
- [153] Tianwei Dai, Sriram Vijayakrishnan, Filip T Szczypiński, Jean-François Ayme, Ehsan Simaei, Thomas Fellowes, Rob Clowes, Lyubomir Kotopantov, Caitlin E Shields, Zhengxue Zhou, et al. Autonomous mobile robots for exploratory synthetic chemistry. *Nature*, pages 1–8, Nov 2024.
- [154] Yahao Dai, Henry Chan, Aikaterini Vriza, Fredrick Kim, Yunfei Wang, Wei Liu, Naisong Shan, Jing Xu, Max Weires, Yukun Wu, et al. Adaptive ai decision interface for autonomous electronic material discovery. *arXiv preprint arXiv:2504.13344*, 2025.
- [155] Preetam Prabhu Srikanth Dammu, Himanshu Naidu, Mouly Dewan, YoungMin Kim, Tanya Roosta, Aman Chadha, and Chirag Shah. Claimver: Explainable claim-level verification and evidence attribution of text through knowledge graphs. *arXiv preprint arXiv:2403.09724*, 2024.
- [156] Raghav Dangayach, Nohyeong Jeong, Elif Demirel, Nigmet Uzal, Victor Fung, and Yongsheng Chen. Machine learning-aided inverse design and discovery of novel polymeric materials for membrane separation. *Environmental Science & Technology*, 59(2):993–1012, Dec 2024.
- [157] Mike D’Arcy, Tom Hope, Larry Birnbaum, and Doug Downey. Marg: Multi-agent review generation for scientific papers. *arXiv preprint arXiv:2401.04259*, 2024.
- [158] Maxime Darrin, Ines Arous, Pablo Piantanida, and Jackie CK Cheung. Glimpse: Pragmatically informative multi-document summarization for scholarly reviews. *arXiv preprint arXiv:2406.07359*, 2024.
- [159] Kourosh Darvish, Marta Skreta, Yuchi Zhao, Naruki Yoshikawa, Sagnik Som, Miroslav Bogdanovic, Yang Cao, Han Hao, Haoping Xu, Alán Aspuru-Guzik, et al. Organa: a robotic assistant for automated chemistry experimentation and characterization. *Matter*, 8(2), Feb 2025.
- [160] Anubrata Das, Houjiang Liu, Venelin Kovatchev, and Matthew Lease. The state of human-centered nlp technology for fact-checking. *Information processing & management*, 60(2):103219, Mar 2023.

- [161] Debajyoti Dasgupta, Arijit Mondal, and Partha Pratim Chakrabarti. Empowering ai as autonomous researchers: Evaluating llms in generating novel research ideas through automated metrics. In *2nd AI4Research Workshop: Towards a Knowledge-grounded Scientific Research Lifecycle*, Dec 2024.
- [162] Filipe de Avila Belbute-Peres, Kevin Smith, Kelsey Allen, Josh Tenenbaum, and J Zico Kolter. End-to-end differentiable physics for learning and control. *Advances in Neural Information Processing Systems*, 31, Dec 2018.
- [163] Nicola De Cao and Thomas Kipf. Molgan: An implicit generative model for small molecular graphs. *arXiv preprint arXiv:1805.11973*, 2018.
- [164] Ties de Kok. Chatgpt for textual analysis? how to use generative llms in accounting research. *Management Science*, Jan 2025. doi: 10.1287/mnsc.2023.03253. URL <https://doi.org/10.1287/mnsc.2023.03253>. Published online in Articles in Advance, 13 Jan 2025.
- [165] Michael S. Deiner, Vlad Honcharov, Jiawei Li, Tim K. Mackey, Travis C. Porco, and Urmimala Sarkar. Large language models can enable inductive thematic analysis of a social media corpus in a single prompt: Human validation study. *JMIR Infodemiology*, 4:e59641, August 2024. doi: 10.2196/59641. URL <https://doi.org/10.2196/59641>.
- [166] Zekun Deng, Zixin Zeng, Weiye Gu, Jiawen Ji, and Bolin Hua. Automatic related work section generation by sentence extraction and reordering. In *AII@ iConference*, pages 101–110, Jan 2021.
- [167] Saaketh Desai, Sadhvikas Addamane, Jeffrey Y Tsao, Igal Brener, Laura P Swiler, Remi Dingreville, and Prasad P Iyer. Autosclab: A self-driving laboratory for interpretable scientific discovery. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 146–154, Apr 2025.
- [168] Jay DeYoung, Iz Beltagy, Madeleine van Zuylen, Bailey Kuehl, and Lucy Lu Wang. Ms2: Multi-document summarization of medical studies. *arXiv preprint arXiv:2104.06486*, 2021.
- [169] Oscar Díaz, Xabier Garmendia, and Juanan Pereira. Streamlining the review process: Ai-generated annotations in research manuscripts. *arXiv preprint arXiv:2412.00281*, 2024.
- [170] Danica Dillion, Niket Tandon, Yuling Gu, and Kurt Gray. Can ai language models replace human participants? *Trends in Cognitive Sciences*, 27(7):597–600, Jul 2023. ISSN 1364-6613. doi: <https://doi.org/10.1016/j.tics.2023.04.008>. URL <https://www.sciencedirect.com/science/article/pii/S1364661323000980>.
- [171] Ning Ding, Shang Qu, Linhai Xie, Yifei Li, Zaoqu Liu, Kaiyan Zhang, Yibai Xiong, Yuxin Zuo, Zhangren Chen, Ermo Hua, et al. Automating exploratory proteomics research via language models. *arXiv preprint arXiv:2411.03743*, 2024.
- [172] Ying Ding and Blaise Cronin. Popular and/or prestigious? measures of scholarly esteem. *Information processing & management*, 47(1):80–96, Jan 2011.
- [173] Martin Docekal, Martin Fajcik, and Pavel Smrz. Oarelatedwork: A large-scale dataset of related work sections with full-texts from open access sources. *arXiv preprint arXiv:2405.01930*, 2024.
- [174] Schmidt Dominik, Jiang Zhengyao, and Wu Yuxiang. Aide: Human-level performance on data science competitions, Apr 2023. URL <https://www.weco.ai/blog/technical-report>. AIDE.

- [175] Tommaso Dorigo, Andrea Giammanco, Pietro Vischia, Max Aehle, Mateusz Bawaj, Alexey Boldyrev, Pablo de Castro Manzano, Denis Derkach, Julien Donini, Auralee Edelen, et al. Toward the end-to-end optimization of particle physics instruments with differentiable programming. *Reviews in Physics*, 10: 100085, Jun 2023.
- [176] Bohdana Daskaliuk, Olena Zimba, Marlen Yessirkepov, Iryna Klishch, and Roman Yatsyshyn. Artificial intelligence in peer review: enhancing efficiency while preserving integrity. *Journal of Korean medical science*, 40(7), Feb 2025.
- [177] Ryan S Dove, Roy J Hartfield, and Mark Carpenter. Semi-supervised classification with novelty detection using support vector machines and linear discriminant analysis. In *AIAA SCITECH 2025 Forum*, page 0705, Jan 2025.
- [178] Mingxuan Du, Benfeng Xu, Chiwei Zhu, Xiaorui Wang, and Zhendong Mao. Deepresearch bench: A comprehensive benchmark for deep research agents. *arXiv preprint arXiv:2506.11763*, 2025.
- [179] Wanyu Du, Zae Myung Kim, Vipul Raheja, Dhruv Kumar, and Dongyeop Kang. Read, revise, repeat: A system demonstration for human-in-the-loop iterative text revision. *arXiv preprint arXiv:2204.03685*, 2022.
- [180] Nils Dycke, Ilia Kuznetsov, and Iryna Gurevych. NLPeer: A unified resource for the computational study of peer review. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5049–5073, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.277. URL <https://aclanthology.org/2023.acl-long.277/>.
- [181] Sana Ebrahimi, Kaiwen Chen, Abolfazl Asudeh, Gautam Das, and Nick Koudas. Axolotl: fairness through assisted self-debiasing of large language model outputs. *arXiv preprint arXiv:2403.00198*, 2024.
- [182] Brice Edelman and Jeffrey Skolnick. Valsci: an open-source, self-hostable literature review utility for automated large-batch scientific claim verification using large language models. *BMC bioinformatics*, 26(1):1–25, May 2025.
- [183] Kristina Edfeldt, Aled M Edwards, Ola Engkvist, Judith Günther, Matthew Hartley, David G Hulcoop, Andrew R Leach, Brian D Marsden, Amelie Menge, Leonie Misquitta, et al. A data science roadmap for open science organizations engaged in early-stage drug discovery. *Nature Communications*, 15(1): 5640, Jul 2024.
- [184] Steffen Eger, Yong Cao, Jennifer D’Souza, Andreas Geiger, Christian Greisinger, Stephanie Gross, Yufang Hou, Brigitte Krenn, Anne Lauscher, Yizhi Li, et al. Transforming science with large language models: A survey on ai-assisted scientific discovery, experimentation, content generation, and evaluation. *arXiv preprint arXiv:2502.05151*, 2025.
- [185] Upol Ehsan and Mark Riedl. Explainable ai reloaded: Challenging the xai status quo in the era of large language models. In *Proceedings of the Halfway to the Future Symposium*, pages 1–8, Oct 2024.
- [186] Christelle Ekosso, Hao Liu, Avery Glagovich, Dustin Nguyen, Sarah Maurer, and Joshua Schrier. Accelerating the discovery of abiotic vesicles with ai-guided automated experimentation. *Langmuir*, 41(1):858–867, 2024.

- [187] Islam Eldifrawi, Shengrui Wang, and Amine Trabelsi. Automated justification production for claim veracity in fact checking: A survey on architectures and approaches. *arXiv preprint arXiv:2407.12853*, 2024.
- [188] Felix Faltings, Michel Galley, Gerold Hintz, Chris Brockett, Chris Quirk, Jianfeng Gao, and Bill Dolan. Text editing by command. *arXiv preprint arXiv:2010.12826*, 2020.
- [189] Angela Fan and Claire Gardent. Generating full length wikipedia biographies: The impact of gender bias on the retrieval-based generation of women biographies. *arXiv preprint arXiv:2204.05879*, 2022.
- [190] Angela Fan, Beliz Gokkaya, Mark Harman, Mitya Lyubarskiy, Shubho Sengupta, Shin Yoo, and Jie M Zhang. Large language models for software engineering: Survey and open problems. In *2023 IEEE/ACM International Conference on Software Engineering: Future of Software Engineering (ICSE-FoSE)*, pages 31–53. IEEE, Oct 2023.
- [191] Wenqi Fan, Yujuan Ding, Liangbo Ning, Shijie Wang, Hengyun Li, Dawei Yin, Tat-Seng Chua, and Qing Li. A survey on rag meeting llms: Towards retrieval-augmented large language models. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 6491–6501, Aug 2024.
- [192] Zhihao Fan, Jialong Tang, Wei Chen, Siyuan Wang, Zhongyu Wei, Jun Xi, Fei Huang, and Jingren Zhou. Ai hospital: Benchmarking large language models in a multi-agent medical interaction simulator. *arXiv preprint arXiv:2402.09742*, 2024.
- [193] You-Le Fang, Dong-Shan Jian, Xiang Li, and Yan-Qing Ma. Ai-newton: A concept-driven physical law discovery system without prior physical knowledge. *arXiv preprint arXiv:2504.01538*, 2025.
- [194] Shai Farber. Enhancing peer review efficiency: A mixed-methods analysis of artificial intelligence-assisted reviewer selection across academic disciplines. *Learned Publishing*, 37(4):e1638, Oct 2024.
- [195] Shai Farber. Enhancing academic decision-making: A pilot study of ai-supported journal selection in higher education. *Innovative Higher Education*, pages 1–19, Feb 2025.
- [196] Manaal Faruqui, Ellie Pavlick, Ian Tenney, and Dipanjan Das. Wikiatomicedits: A multilingual corpus of wikipedia edits for modeling language and discourse. *arXiv preprint arXiv:1808.09422*, 2018.
- [197] Yao Fehlis. Uncovering bottlenecks and optimizing scientific lab workflows with cycle time reduction agents. *arXiv preprint arXiv:2505.21534*, 2025.
- [198] Yao Fehlis, Paul Mandel, Charles Crain, Betty Liu, and David Fuller. Accelerating drug discovery with artificial: a whole-lab orchestration and scheduling system for self-driving labs. *arXiv preprint arXiv:2504.00986*, 2025.
- [199] Bin Feng, Zequn Liu, Nanlan Huang, Zhiping Xiao, Haomiao Zhang, Srбуhi Mirzoyan, Hanwen Xu, Jiaran Hao, Yinghui Xu, Ming Zhang, et al. A bioactivity foundation model using pairwise meta-learning. *Nature Machine Intelligence*, 6(8):962–974, Aug 2024.
- [200] Jingsen Feng, Ran Xu, and Xu Chu. Openfoamgpt 2.0: end-to-end, trustworthy automation for computational fluid dynamics. *arXiv preprint arXiv:2504.19338*, 2025.
- [201] Kehua Feng, Keyan Ding, Weijie Wang, Xiang Zhuang, Zeyuan Wang, Ming Qin, Yu Zhao, Jianhua Yao, Qiang Zhang, and Huajun Chen. Sciknoweval: Evaluating multi-level scientific knowledge of large language models. *arXiv preprint arXiv:2406.09098*, 2024.

- [202] KJ Feng, Kevin Pu, Matt Latzke, Tal August, Pao Siangliulue, Jonathan Bragg, Daniel S Weld, Amy X Zhang, and Joseph Chee Chang. Cocoa: Co-planning and co-execution with ai agents. *arXiv preprint arXiv:2412.10999*, 2024.
- [203] Ruozhu Feng, Yangang Liang, Tianzhixi Yin, Peiyuan Gao, and Wei Wang. Agentic assistant for material scientists. Apr 2025.
- [204] Tao Feng, Yihang Sun, and Jiaxuan You. Grapheval: A lightweight graph-based llm framework for idea evaluation. *arXiv preprint arXiv:2503.12600*, 2025.
- [205] Leandro Carísio Fernandes, Gustavo Bartz Guedes, Thiago Soares Laitz, Thales Sales Almeida, Rodrigo Nogueira, Roberto Lotufo, and Jayr Pereira. Surveysum: A dataset for summarizing multiple scientific articles into a survey section. In *Brazilian Conference on Intelligent Systems*, pages 431–444. Springer, Jan 2024.
- [206] Emilio Ferrara. Fairness and bias in artificial intelligence: A brief survey of sources, impacts, and mitigation strategies. *Sci*, 6(1):3, Dec 2023.
- [207] Financial Times. Deepmind and biontech build ai lab assistants for scientific research. *Financial Times*, Oct 2024. URL <https://www.ft.com/content/64b1bb33-095e-4cc5-a911-50df76fa3d1d>.
- [208] Emily First, Markus N Rabe, Talia Ringer, and Yuriy Brun. Baldur: Whole-proof generation and repair with large language models. In *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, pages 1229–1241, Nov 2023.
- [209] Marcio Fonseca and Shay Cohen. Can large language model summarizers adapt to diverse scientific communication goals? In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Findings of the Association for Computational Linguistics: ACL 2024*, pages 8599–8618, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.508. URL <https://aclanthology.org/2024.findings-acl.508/>.
- [210] Thibault Formal, Carlos Lassance, Benjamin Piwowarski, and Stéphane Clinchant. Splade v2: Sparse lexical and expansion model for information retrieval. *arXiv preprint arXiv:2109.10086*, 2021.
- [211] Maria Frasca, Davide La Torre, Gabriella Pravettoni, and Ilaria Cutica. Explainable and interpretable artificial intelligence in medicine: a systematic bibliometric review. *Discover Artificial Intelligence*, 4(1):15, Feb 2024.
- [212] Simon Frieder, Jonas Bayer, Katherine M Collins, Julius Berner, Jacob Loader, András Juhász, Fabian Ruehle, Sean Welleck, Gabriel Poesia, Ryan-Rhys Griffiths, et al. Data for mathematical copilots: Better ways of presenting proofs for machine learning. *arXiv preprint arXiv:2412.15184*, 2024.
- [213] Samuel Frontull and Georg Moser. Rule-based, neural and llm back-translation: Comparative insights from a variant of ladin. *arXiv preprint arXiv:2407.08819*, 2024.
- [214] Yongfan Fu, Jian Luo, Guofang Nan, and Dahui Li. Peer review expert group recommendation: A multi-subject coverage-based approach. *Expert Systems with Applications*, 264:125971, Mar 2025.
- [215] Carlo Galli, Chiara Moretti, and Elena Calciolari. Intelligent summaries: Will artificial intelligence mark the finale for biomedical literature reviews? *Learned Publishing*, Dec 2024.

- [216] Shubham Gandhi, Dhruv Shah, Manasi Patwardhan, Lovekesh Vig, and Gautam Shroff. Research-codeagent: An llm multi-agent system for automated codification of research methodologies. *arXiv preprint arXiv:2504.20117*, 2025.
- [217] Amrita Ganguly, Aditya Johri, Areej Ali, and Nora McDonald. Generative artificial intelligence for academic research: evidence from guidance issued for researchers by higher education institutions in the united states. *AI and Ethics*, pages 1–17, Mar 2025.
- [218] Debargha Ganguly, Vikash Singh, Sreehari Sankar, Biyao Zhang, Xuecen Zhang, Srinivasan Iyengar, Xiaotian Han, Amit Sharma, Shivkumar Kalyanaraman, and Vipin Chaudhary. Grammars of formal uncertainty: When to trust llms in automated reasoning tasks. *arXiv preprint arXiv:2505.20047*, 2025.
- [219] Amit Gangwal, Azim Ansari, Iqar Ahmad, Abul Kalam Azad, and Wan Mohd Azizi Wan Sulaiman. Current strategies to address data scarcity in artificial intelligence-based drug discovery: A comprehensive review. *Computers in Biology and Medicine*, 179:108734, Sep 2024.
- [220] Xian Gao, Jiacheng Ruan, Jingsheng Gao, Ting Liu, and Yuzhuo Fu. Reviewagents: Bridging the gap between human and ai-generated paper reviews. *arXiv preprint arXiv:2503.08506*, 2025.
- [221] Xian Gao, Zongyun Zhang, Mingye Xie, Ting Liu, and Yuzhuo Fu. Graph of ai ideas: Leveraging knowledge graphs and llms for ai research idea generation. *arXiv preprint arXiv:2503.08549*, 2025.
- [222] Andres Garcia-Silva, Cristian Berrio, Jose Manuel Gomez-Perez, Jose Antonio Martínez-Heras, Alessandro Donati, and Ilaria Roma. Spaceqa: Answering questions about the design of space missions and space craft concepts. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 3306–3311, Jul 2022.
- [223] Daniel Garijo, Qifan Yang, Hernán Vargas, Shruti P Gadewar, Kevin Low, Varun Ratnakar, Maximiliano Osorio, Alyssa H Zhu, Agnes McMahon, Yolanda Gil, et al. Neurodisk: An ai approach to automate continuous inquiry-driven discoveries in neuroimaging genetics. *bioRxiv*, Feb 2025.
- [224] Aniketh Garikaparathi, Manasi Patwardhan, Aditya Sanjiv Kanade, Aman Hassan, Lovekesh Vig, and Arman Cohan. Mir: Methodology inspiration retrieval for scientific research problems. *arXiv preprint arXiv:2506.00249*, 2025.
- [225] Aniketh Garikaparathi, Manasi Patwardhan, Lovekesh Vig, and Arman Cohan. Iris: Interactive research ideation system for accelerating scientific discovery. *arXiv preprint arXiv:2504.16728*, 2025.
- [226] Yubin Ge, Ly Dinh, Xiaofeng Liu, Jinsong Su, Ziyao Lu, Ante Wang, and Jana Diesner. Baco: A background knowledge-and content-based framework for citing sentence generation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1466–1478, Aug 2021.
- [227] Mingmeng Geng and Roberto Trotta. Human-llm coevolution: Evidence from academic writing. *arXiv preprint arXiv:2502.09606*, 2025.
- [228] Elizabeth Oommen George. How paperpal enhances english writing quality and improves productivity for japanese academics. Aug 2024.

- [229] Katy Ilonka Gero, Vivian Liu, and Lydia Chilton. Sparks: Inspiration for science writing using language models. In *Proceedings of the 2022 ACM Designing Interactive Systems Conference*, pages 1002–1019, Jun 2022.
- [230] Alireza Ghafarollahi and Markus J Buehler. Protagents: protein discovery via large language model multi-agent collaborations combining physics and machine learning. *Digital Discovery*, 3(7):1389–1409, May 2024.
- [231] Alireza Ghafarollahi and Markus J Buehler. Sciagents: Automating scientific discovery through multi-agent intelligent graph reasoning. *arXiv preprint arXiv:2409.05556*, 2024.
- [232] Alireza Ghafarollahi and Markus J Buehler. Sparks: Multi-agent artificial intelligence model discovers protein design principles. *arXiv preprint arXiv:2504.19017*, 2025.
- [233] Ali Essam Ghareeb, Benjamin Chang, Ludovico Mitchener, Angela Yiu, Caralyn J Szostkiewicz, Jon M Laurent, Muhammed T Razzak, Andrew D White, Michaela M Hinks, and Samuel G Rodriques. Robin: A multi-agent system for automating scientific discovery. *arXiv preprint arXiv:2505.13400*.
- [234] Ali Gharizadeh, Karim Abbasi, Amin Ghareyazi, Mohammad RK Mofrad, and Hamid R Rabiee. Hgtdr: Advancing drug repurposing with heterogeneous graph transformers. *Bioinformatics*, 40(7):btac349, Jul 2024.
- [235] Tirthankar Ghosal, Sandeep Kumar, Prabhat Kumar Bharti, and Asif Ekbal. Peer review analyze: A novel benchmark resource for computational analysis of peer reviews. *Plos one*, 17(1):e0259238, Jan 2022.
- [236] Asish Ghoshal, Srinivasan Iyer, Bhargavi Paranjape, Kushal Lakhotia, Scott Wen-tau Yih, and Yashar Mehdad. Quaser: Question answering with scalable extractive rationalization. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 1208–1218, Jul 2022.
- [237] Karan Girotra, Lennart Meincke, Christian Terwiesch, and Karl T Ulrich. Ideas are dimes a dozen: Large language models for idea generation in innovation. *The Wharton School Research Paper Forthcoming*, Jul 2023.
- [238] Moshe Glickman and Tali Sharot. How human–ai feedback loops alter human perceptual, emotional and social judgements. *Nature Human Behaviour*, 9(2):345–359, 2025.
- [239] Max Glockner, Yufang Hou, and Iryna Gurevych. Missing counter-evidence renders nlp fact-checking unrealistic for misinformation. *arXiv preprint arXiv:2210.13865*, 2022.
- [240] Max Glockner, Yufang Hou, Preslav Nakov, and Iryna Gurevych. Grounding fallacies misrepresenting scientific publications in evidence. *arXiv preprint arXiv:2408.12812*, 2024.
- [241] Ozan Gokdemir, Carlo Siebenschuh, Alexander Brace, Azton Wells, Brian Hsu, Kyle Hippe, Priyanka V Setty, Aswathy Ajith, J Gregory Pauloski, Varuni Sastry, et al. Hiperrag: High-performance retrieval augmented generation for scientific insights. *arXiv preprint arXiv:2505.04846*, Jun 2025.
- [242] Ali Goli and Amandeep Singh. Frontiers: Can large language models capture human preferences? *Marketing Science*, 43(4):709–722, Apr 2024. doi: 10.1287/mksc.2023.0306. URL <https://doi.org/10.1287/mksc.2023.0306>.

- [243] Catalina Gomez, Sue Min Cho, Shichang Ke, Chien-Ming Huang, and Mathias Unberath. Human-ai collaboration is not very collaborative yet: A taxonomy of interaction patterns in ai-assisted decision making from a systematic review. *Frontiers in Computer Science*, 6:1521066, Jan 2025.
- [244] Rafael Gómez-Bombarelli, Jennifer N Wei, David Duvenaud, José Miguel Hernández-Lobato, Benjamín Sánchez-Lengeling, Dennis Sheberla, Jorge Aguilera-Iparraguirre, Timothy D Hirzel, Ryan P Adams, and Alán Aspuru-Guzik. Automatic chemical design using a data-driven continuous representation of molecules. *ACS central science*, 4(2):268–276, Jan 2018.
- [245] Jose Manuel Gomez-Perez and Raul Ortega. Look, read and enrich-learning from scientific figures and their captions. In *Proceedings of the 10th International Conference on Knowledge Capture*, pages 101–108, Sep 2019.
- [246] Rubén González-Sendino, Emilio Serrano, and Javier Bajo. Mitigating bias in artificial intelligence: Fair data generation via causal models for transparent and explainable decision-making. *Future Generation Computer Systems*, 155:384–401, Jun 2024.
- [247] Zachary Goodsell and Juhani Yli-Vakkuri. Lf: a foundational higher-order-logic. *arXiv preprint arXiv:2401.11050*, 2024.
- [248] Diego Gosmar and Deborah A Dahl. Hallucination mitigation using agentic ai natural language-based frameworks. *arXiv preprint arXiv:2501.13946*, 2025.
- [249] 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.
- [250] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- [251] Samuel Greydanus, Misko Dzamba, and Jason Yosinski. Hamiltonian neural networks. *Advances in Neural Information Processing Systems*, 32, Jul 2019.
- [252] 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.
- [253] Antoine Grosnit, Alexandre Maraval, James Doran, Giuseppe Paolo, Albert Thomas, Refinath Shahul Hameed Nabeezath Beevi, Jonas Gonzalez, Khyati Khandelwal, Ignacio Iacobacci, Abdelhakim Benechehab, et al. Large language models orchestrating structured reasoning achieve kaggle grand-master level. *arXiv preprint arXiv:2411.03562*, 2024.
- [254] Ken Gu, Ruoxi Shang, Ruien Jiang, Keying Kuang, Richard-John Lin, Donghe Lyu, Yue Mao, Youran Pan, Teng Wu, Jiaqian Yu, et al. Blade: Benchmarking language model agents for data-driven science. *arXiv preprint arXiv:2408.09667*, 2024.
- [255] Nianlong Gu and Richard HR Hahnloser. Controllable citation sentence generation with language models. *arXiv preprint arXiv:2211.07066*, 2022.
- [256] Tianyang Gu, Jingjin Wang, Zhihao Zhang, and HaoHong Li. Llms can realize combinatorial creativity: generating creative ideas via llms for scientific research. *arXiv preprint arXiv:2412.14141*, 2024.

- [257] Xuemei Gu and Mario Krenn. Generation and human-expert evaluation of interesting research ideas using knowledge graphs and large language models. *arXiv preprint arXiv:2405.17044*, 2024.
- [258] Xuemei Gu and Mario Krenn. Interesting scientific idea generation using knowledge graphs and llms: Evaluations with 100 research group leaders. *arXiv preprint arXiv:2405.17044*, 2024.
- [259] Yuan Guan, Jakkapong Inchai, Zhuoqing Fang, Jacky Law, Alberto Alonzo Garcia Brito, Annalisa Pawlosky, Juraj Gottweis, Alexander Daryin, Artiom Myaskovsky, Anil Palepu, et al. Ai-assisted drug re-purposing for human liver fibrosis. *bioRxiv*, pages 2025–04, May 2025.
- [260] Badra Souhila Guendouzi, Samir Ouchani, Hiba EL Assaad, and Madeleine EL Zaher. A systematic review of federated learning: Challenges, aggregation methods, and development tools. *Journal of Network and Computer Applications*, 220:103714, Nov 2023.
- [261] Maeva Guerrier, Karthik Soma, Hassan Fouad, and Giovanni Beltrame. Guided by guardrails: Control barrier functions as safety instructors for robotic learning. *arXiv preprint arXiv:2505.18858*, 2025.
- [262] Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Yu Wu, YK Li, et al. Deepseek-coder: When the large language model meets programming—the rise of code intelligence. *arXiv preprint arXiv:2401.14196*, 2024.
- [263] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- [264] Siyuan Guo, Cheng Deng, Ying Wen, Hechang Chen, Yi Chang, and Jun Wang. Ds-agent: Automated data science by empowering large language models with case-based reasoning. *arXiv preprint arXiv:2402.17453*, 2024.
- [265] Yue Guo, Wei Qiu, Yizhong Wang, and Trevor Cohen. Automated lay language summarization of biomedical scientific reviews. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 160–168, May 2021.
- [266] Ziyu Guo, Ray Zhang, Hao Chen, Jialin Gao, Dongzhi Jiang, Jiaze Wang, and Pheng-Ann Heng. Sciverse: Unveiling the knowledge comprehension and visual reasoning of llms on multi-modal scientific problems. *arXiv preprint arXiv:2503.10627*, 2025.
- [267] Rohan Gupta, Devesh Srivastava, Mehar Sahu, Swati Tiwari, Rashmi K Ambasta, and Pravir Kumar. Artificial intelligence to deep learning: machine intelligence approach for drug discovery. *Molecular diversity*, 25:1315–1360, Apr 2021.
- [268] Tarun Gupta and Danish Pruthi. All that glitters is not novel: Plagiarism in ai generated research. *arXiv preprint arXiv:2502.16487*, 2025.
- [269] Anthony GX-Chen, Dongyan Lin, Mandana Samiei, Doina Precup, Blake A Richards, Rob Fergus, and Kenneth Marino. Language agents mirror human causal reasoning biases. how can we help them think like scientists? *arXiv preprint arXiv:2505.09614*, 2025.
- [270] Hendrik Haarmann. Enhance innovation by boosting idea generation with large language models. *INFORMS Journal on Computing*, Jul 2025.

- [271] Thilo Hagendorff, Sarah Fabi, and Michal Kosinski. Human-like intuitive behavior and reasoning biases emerged in large language models but disappeared in chatgpt. *Nature Computational Science*, 3(10):833–838, October 2023. ISSN 2662-8457. doi: 10.1038/s43588-023-00527-x. URL <http://dx.doi.org/10.1038/s43588-023-00527-x>.
- [272] Lewis Hammond, Alan Chan, Jesse Clifton, Jason Hoelscher-Obermaier, Akbir Khan, Euan McLean, Chandler Smith, Wolfram Barfuss, Jakob Foerster, Tomáš Gavenčík, et al. Multi-agent risks from advanced ai. *arXiv preprint arXiv:2502.14143*, 2025.
- [273] Viktor Hangya, Hossain Shaikh Saadi, and Alexander Fraser. Improving low-resource languages in pre-trained multilingual language models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 11993–12006, Jan 2022.
- [274] Matthew G Hanna, Liron Pantanowitz, Brian Jackson, Octavia Palmer, Shyam Visweswaran, Joshua Pantanowitz, Mustafa Deebajah, and Hooman H Rashidi. Ethical and bias considerations in artificial intelligence/machine learning. *Modern Pathology*, 38(3):100686, Mar 2025.
- [275] Qian Yue Hao, Jingyang Fan, Fengli Xu, Jian Yuan, and Yong Li. Hlm-cite: Hybrid language model workflow for text-based scientific citation prediction. *arXiv preprint arXiv:2410.09112*, 2024.
- [276] Shibo Hao, Tianyang Liu, Zhen Wang, and Zhiting Hu. Toolkengpt: Augmenting frozen language models with massive tools via tool embeddings. *Advances in Neural Information Processing Systems*, 36:45870–45894, Dec 2023.
- [277] Kenneth D Harris. Airus: a simple workflow for ai-assisted exploration of scientific data. *bioRxiv*, pages 2025–02, Feb 2025.
- [278] Russell Hartley. Efficacy analysis of online artificial intelligence fact-checking tools. *The International Review of Information Ethics*, 33(1), Apr 2024.
- [279] Helia Hashemi, Jason Eisner, Corby Rosset, Benjamin Van Durme, and Chris Kedzie. LLM-rubric: A multidimensional, calibrated approach to automated evaluation of natural language texts. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13806–13834, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.745. URL <https://aclanthology.org/2024.acl-long.745/>.
- [280] Vikas Hassija, Vinay Chamola, Atmesh Mahapatra, Abhinandan Singal, Divyansh Goel, Kaizhu Huang, Simone Scardapane, Indro Spinelli, Mufti Mahmud, and Amir Hussain. Interpreting black-box models: a review on explainable artificial intelligence. *Cognitive Computation*, 16(1):45–74, Aug 2024.
- [281] Kan Hatakeyama-Sato, Toshihiko Nishida, Kenta Kitamura, Yoshitaka Ushiku, Koichi Takahashi, Yuta Nabae, and Teruaki Hayakawa. Perspective on utilizing foundation models for laboratory automation in materials research. *arXiv preprint arXiv:2506.12312*, 2025.
- [282] Jesse Haworth, Rishi Biswas, Justin Opfermann, Michael Kam, Yaning Wang, Desire Pantalone, Francis X Creighton, Robin Yang, Jin U Kang, and Axel Krieger. Autonomous robotic system with optical coherence tomography guidance for vascular anastomosis. *arXiv preprint arXiv:2410.07493*, 2024.

- [283] Thomas Hayes, Roshan Rao, Halil Akin, Nicholas J Sofroniew, Deniz Oktay, Zeming Lin, Robert Verkuil, Vincent Q Tran, Jonathan Deaton, Marius Wiggert, et al. Simulating 500 million years of evolution with a language model. *Science*, page eads0018, Jan 2025.
- [284] Kaiyu He and Zhiyu Chen. From reasoning to learning: A survey on hypothesis discovery and rule learning with large language models. *arXiv preprint arXiv:2505.21935*, 2025.
- [285] Yichen He, Guanhua Huang, Peiyuan Feng, Yuan Lin, Yuchen Zhang, Hang Li, et al. Pasa: An llm agent for comprehensive academic paper search. *arXiv preprint arXiv:2501.10120*, 2025.
- [286] Tina Heger, Alsayed Algergawy, Marc Brinner, Jonathan M Jeschke, Birgitta König-Ries, Daniel Mietchen, and Sina Zarriß. Natural language hypotheses in scientific papers and how to tame them: Suggested steps for formalizing complex scientific claims. In *Conference on Advances in Robust Argumentation Machines*, pages 3–19. Springer Nature Switzerland Cham, Jun 2024.
- [287] Michael V Heinz, Daniel M Mackin, Brianna M Trudeau, Sukanya Bhattacharya, Yinzhou Wang, Haley A Banta, Abi D Jewett, Abigail J Salzhauer, Tess Z Griffin, and Nicholas C Jacobson. Randomized trial of a generative ai chatbot for mental health treatment. *Nejm Ai*, 2(4):Aloa2400802, Mar 2025.
- [288] Sam Henry and Bridget T McInnes. Literature based discovery: models, methods, and trends. *Journal of biomedical informatics*, 74:20–32, Oct 2017.
- [289] Emily Herron, Junqi Yin, and Feiyi Wang. Scitrust: Evaluating the trustworthiness of large language models for science. In *SC24-W: Workshops of the International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 72–78. IEEE, Nov 2024.
- [290] Lukas Hilgert, Danni Liu, and Jan Niehues. Evaluating and training long-context large language models for question answering on scientific papers. In Sachin Kumar, Vidhisha Balachandran, Chan Young Park, Weijia Shi, Shirley Anugrah Hayati, Yulia Tsvetkov, Noah Smith, Hannaneh Hajishirzi, Dongyeop Kang, and David Jurgens, editors, *Proceedings of the 1st Workshop on Customizable NLP: Progress and Challenges in Customizing NLP for a Domain, Application, Group, or Individual (CustomNLP4U)*, pages 220–236, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.customnlp4u-1.17. URL <https://aclanthology.org/2024.customnlp4u-1.17/>.
- [291] Cong Duy Vu Hoang and Min-Yen Kan. Towards automated related work summarization. In *Coling 2010: Posters*, pages 427–435, Aug 2010.
- [292] Michael Peter Hoffmann, Jan Fillies, and Adrian Paschke. Malinowski in the age of ai: Can large language models create a text game based on an anthropological classic? *arXiv preprint arXiv:2410.20536*, 2024.
- [293] Brendan Hogan, Anmol Kabra, Felipe Siqueira Pacheco, Laura Greenstreet, Joshua Fan, Aaron Ferber, Marta Ummus, Alecsander Brito, Olivia Graham, Lillian Aoki, et al. Aiscivision: A framework for specializing large multimodal models in scientific image classification. *arXiv preprint arXiv:2410.21480*, 2024.
- [294] Steffen Holter and Mennatallah El-Assady. Deconstructing human-ai collaboration: Agency, interaction, and adaptation. In *Computer Graphics forum*, volume 43, page e15107. Wiley Online Library, Jun 2024.
- [295] Niklas Höpner, Leon Eshuijs, Dimitrios Alivanistos, Giacomo Zamprogno, and Ilaria Tiddi. Automatic evaluation metrics for artificially generated scientific research. *arXiv preprint arXiv:2503.05712*, 2025.

- [296] Eftekhari Hossain, Sanjeev Kumar Sinha, Naman Bansal, R. Alexander Knipper, Souvika Sarkar, John Salvador, Yash Mahajan, Sri Ram Pavan Kumar Guttikonda, Mousumi Akter, Md. Mahadi Hassan, Matthew Freestone, Matthew C. Williams Jr., Dongji Feng, and Santu Karmaker. LLMs as meta-reviewers' assistants: A case study. In Luis Chiruzzo, Alan Ritter, and Lu Wang, editors, *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 7763–7803, Albuquerque, New Mexico, April 2025. Association for Computational Linguistics. ISBN 979-8-89176-189-6. doi: 10.18653/v1/2025.naacl-long.395. URL <https://aclanthology.org/2025.naacl-long.395/>.
- [297] Soodeh Hosseini and Hossein Seilani. The role of agentic ai in shaping a smart future: A systematic review. *Array*, page 100399, Jul 2025.
- [298] Jhih-Yi Hsieh. *Automated Peer-Reviewer Assignment can be Manipulated to Secure Reviews from Colluders*. PhD thesis, Carnegie Mellon University Pittsburgh, PA, May 2024.
- [299] Chao-Chun Hsu, Erin Bransom, Jenna Sparks, Bailey Kuehl, Chenhao Tan, David Wadden, Lucy Lu Wang, and Aakanksha Naik. Chime: Llm-assisted hierarchical organization of scientific studies for literature review support. *arXiv preprint arXiv:2407.16148*, 2024.
- [300] Ting-Yao Hsu, Chieh-Yang Huang, Shih-Hong Huang, Ryan Rossi, Sungchul Kim, Tong Yu, C Lee Giles, and Ting-Hao Kenneth Huang. Scicapenter: Supporting caption composition for scientific figures with machine-generated captions and ratings. In *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems*, pages 1–9, May 2024.
- [301] Linmei Hu, Duokang Wang, Yiming Pan, Jifan Yu, Yingxia Shao, Chong Feng, and Liqiang Nie. Novachart: A large-scale dataset towards chart understanding and generation of multimodal large language models. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 3917–3925, Oct 2024. URL <https://openreview.net/forum?id=PTYL6011vp>.
- [302] Mengkang Hu, Yao Mu, Xinmiao Yu, Mingyu Ding, Shiguang Wu, Wenqi Shao, Qiguang Chen, Bin Wang, Yu Qiao, and Ping Luo. Tree-planner: Efficient close-loop task planning with large language models. *arXiv preprint arXiv:2310.08582*, 2023.
- [303] Mengkang Hu, Tianxing Chen, Qiguang Chen, Yao Mu, Wenqi Shao, and Ping Luo. Hiagent: Hierarchical working memory management for solving long-horizon agent tasks with large language model. *arXiv preprint arXiv:2408.09559*, 2024.
- [304] Mengkang Hu, Tianxing Chen, Yude Zou, Yuheng Lei, Qiguang Chen, Ming Li, Yao Mu, Hongyuan Zhang, Wenqi Shao, and Ping Luo. Text2world: Benchmarking large language models for symbolic world model generation. *arXiv preprint arXiv:2502.13092*, 2025.
- [305] Mengkang Hu, Yuhang Zhou, Wendong Fan, Yuzhou Nie, Bowei Xia, Tao Sun, Ziyu Ye, Zhaoxuan Jin, Yingru Li, Qiguang Chen, et al. Owl: Optimized workforce learning for general multi-agent assistance in real-world task automation. *arXiv preprint arXiv:2505.23885*, 2025.
- [306] Xiang Hu, Hongyu Fu, Jing Wang, Yifeng Wang, Zhikun Li, Renjun Xu, Yu Lu, Yaochu Jin, Lili Pan, and Zhenzhong Lan. Nova: An iterative planning and search approach to enhance novelty and diversity of llm generated ideas. *arXiv preprint arXiv:2410.14255*, 2024.
- [307] Xueyu Hu, Ziyu Zhao, Shuang Wei, Ziwei Chai, Qianli Ma, Guoyin Wang, Xuwu Wang, Jing Su, Jingjing Xu, Ming Zhu, et al. Infiagent-dabench: Evaluating agents on data analysis tasks. *arXiv preprint arXiv:2401.05507*, 2024.

- [308] Yue Hu and Xiaojun Wan. Automatic generation of related work sections in scientific papers: an optimization approach. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1624–1633, Oct 2014.
- [309] Yuntong Hu, Zhuofeng Li, Zheng Zhang, Chen Ling, Raasikh Kanjiani, Boxin Zhao, and Liang Zhao. Taxonomy tree generation from citation graph. *arXiv preprint arXiv:2410.03761*, 2024.
- [310] Tianyu Hua, Harper Hua, Violet Xiang, Benjamin Klieger, Sang T Truong, Weixin Liang, Fan-Yun Sun, and Nick Haber. Researchcodebench: Benchmarking llms on implementing novel machine learning research code. *arXiv preprint arXiv:2506.02314*, 2025.
- [311] Jiaxin Huang, Shixiang Shane Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. Large language models can self-improve. *arXiv preprint arXiv:2210.11610*, 2022.
- [312] Jincai Huang, Yongjun Xu, Qi Wang, Qi Cheems Wang, Xingxing Liang, Fei Wang, Zhao Zhang, Wei Wei, Boxuan Zhang, Libo Huang, et al. Foundation models and intelligent decision-making: Progress, challenges, and perspectives. *The Innovation*, Jun 2025.
- [313] Jinyang Huang, Xiachong Feng, Qiguang Chen, Hanjie Zhao, Zihui Cheng, Jiesong Bai, Jingxuan Zhou, Min Li, and Libo Qin. Mldebugging: Towards benchmarking code debugging across multi-library scenarios. *arXiv preprint arXiv:2506.13824*, 2025.
- [314] 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.
- [315] Kexin Huang, Payal Chandak, Qianwen Wang, Shreyas Havaladar, Akhil Vaid, Jure Leskovec, Girish N Nadkarni, Benjamin S Glicksberg, Nils Gehlenborg, and Marinka Zitnik. A foundation model for clinician-centered drug repurposing. *Nature Medicine*, 30(12):3601–3613, Sep 2024.
- [316] Kexin Huang, Serena Zhang, Hanchen Wang, Yuanhao Qu, Yingzhou Lu, Yusuf Roohani, Ryan Li, Lin Qiu, Junze Zhang, Yin Di, et al. Biomni: A general-purpose biomedical ai agent. *bioRxiv*, pages 2025–05, Jun 2025.
- [317] Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *ACM Transactions on Information Systems*, 43(2):1–55, Jan 2025.
- [318] Lifu Huang, Danai Koutra, Adithya Kulkarni, Temiloluwa Prioleau, Qingyun Wu, Yujun Yan, Yaoqing Yang, James Zou, and Dawei Zhou. Towards agentic ai for science: Hypothesis generation, comprehension, quantification, and validation. In *Companion Proceedings of the ACM on Web Conference 2025*, pages 1639–1642, May 2025.
- [319] Muye Huang, Lingling Zhang, Jie Ma, Han Lai, Fangzhi Xu, Yifei Li, Wenjun Wu, Yaqiang Wu, and Jun Liu. Chatsketcher: Reasoning with multimodal feedback and reflection for chart understanding. *arXiv preprint arXiv:2505.19076*, 2025.
- [320] Qian Huang, Jian Vora, Percy Liang, and Jure Leskovec. Mlagentbench: Evaluating language agents on machine learning experimentation. *arXiv preprint arXiv:2310.03302*, 2023.

- [321] Shengzhi Huang, Qicong Wang, Wei Lu, Lingyu Liu, Zhenzhen Xu, and Yong Huang. Papereval: A universal, quantitative, and explainable paper evaluation method powered by a multi-agent system. *Information Processing & Management*, 62(6):104225, Nov 2025.
- [322] Xiang Huang, CY Zhao, Hong Wang, and Shenghong Ju. Ai-assisted inverse design of sequence-ordered high intrinsic thermal conductivity polymers. *Materials Today Physics*, 44:101438, May 2024.
- [323] Yinya Huang, Xiaohan Lin, Zhengying Liu, Qingxing Cao, Huajian Xin, Haiming Wang, Zhenguo Li, Linqi Song, and Xiaodan Liang. Mustard: Mastering uniform synthesis of theorem and proof data. *arXiv preprint arXiv:2402.08957*, 2024.
- [324] Yizhe Huang, Xingbo Wang, Hao Liu, Fanqi Kong, Aoyang Qin, Min Tang, Xiaoxi Wang, Song-Chun Zhu, Mingjie Bi, Siyuan Qi, et al. Adasociety: An adaptive environment with social structures for multi-agent decision-making. *arXiv preprint arXiv:2411.03865*, 2024.
- [325] B Ian Hutchins, Xin Yuan, James M Anderson, and George M Santangelo. Relative citation ratio (rcr): a new metric that uses citation rates to measure influence at the article level. *PLoS biology*, 14(9): e1002541, Sep 2016.
- [326] Holland Hysmith, Elham Foadian, Shakti P Padhy, Sergei V Kalinin, Rob G Moore, Olga S Ovchinnikova, and Mahshid Ahmadi. The future of self-driving laboratories: from human in the loop interactive ai to gamification. *Digital Discovery*, 3(4):621–636, Mar 2024.
- [327] Maximilian Idahl and Zahra Ahmadi. Openreviewer: A specialized large language model for generating critical scientific paper reviews. *arXiv preprint arXiv:2412.11948*, 2024.
- [328] Tal Ifargan, Lukas Hafner, Maor Kern, Ori Alcalay, and Roy Kishony. Autonomous llm-driven research—from data to human-verifiable research papers. *NEJM AI*, 2(1):AIoa2400555, Dec 2025.
- [329] Hayato Ikoma and Teruko Mitamura. Can ai examine novelty of patents?: Novelty evaluation based on the correspondence between patent claim and prior art. *arXiv preprint arXiv:2502.06316*, 2025.
- [330] Autoscience Institute. Carl technical report, Mar 2025. URL <https://drive.google.com/file/d/1iVedOdZDuEdjs41cm9Z7i8oEDGWfzVJq/view>. Carl Technical Report.
- [331] Amazon Artificial General Intelligence. The amazon nova family of models: Technical report and model card. Jun 2025.
- [332] Mete Ismayilzada, Debjit Paul, Antoine Bosselut, and Lonneke van der Plas. Creativity in ai: Progresses and challenges. *arXiv preprint arXiv:2410.17218*, 2024.
- [333] Takumi Ito, Tatsuki Kuribayashi, Hayato Kobayashi, Ana Brassard, Masato Hagiwara, Jun Suzuki, and Kentaro Inui. Diamonds in the rough: Generating fluent sentences from early-stage drafts for academic writing assistance. *arXiv preprint arXiv:1910.09180*, 2019.
- [334] Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Hel- yar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. *arXiv preprint arXiv:2412.16720*, 2024.
- [335] Nazanin Jafari and James Allan. Robust claim verification through fact detection. *arXiv preprint arXiv:2407.18367*, 2024.

- [336] Peter Jansen, Marc-Alexandre Côté, Tushar Khot, Erin Bransom, Bhavana Dalvi Mishra, Bodhisattwa Prasad Majumder, Oyvind Tafjord, and Peter Clark. Discoveryworld: A virtual environment for developing and evaluating automated scientific discovery agents. *Advances in Neural Information Processing Systems*, 37:10088–10116, Dec 2024.
- [337] Rasika Jayarathna, Thossaporn Onsree, Samuel Drummond, Jennifer Naglic, and Jochen Lauterbach. Experimental discovery of novel ammonia synthesis catalysts via active learning. *Journal of Materials Chemistry A*, 12(5):3046–3060, Feb 2024.
- [338] Deyi Ji, Lanyun Zhu, Siqi Gao, Peng Xu, Hongtao Lu, Jieping Ye, and Feng Zhao. Tree-of-table: Unleashing the power of llms for enhanced large-scale table understanding. *arXiv preprint arXiv:2411.08516*, 2024.
- [339] Qinjin Jia, Jialin Cui, Yunkai Xiao, Chengyuan Liu, Parvez Rashid, and Edward Gehring. All-in-one: Multi-task learning bert models for evaluating peer assessments. *International Educational Data Mining Society*, Oct 2021.
- [340] Albert Q Jiang, Sean Welleck, Jin Peng Zhou, Wenda Li, Jiacheng Liu, Mateja Jamnik, Timothée Lacroix, Yuhuai Wu, and Guillaume Lample. Draft, sketch, and prove: Guiding formal theorem provers with informal proofs. *arXiv preprint arXiv:2210.12283*, 2022.
- [341] Albert Qiaochu Jiang, Wenda Li, Szymon Tworkowski, Konrad Czechowski, Tomasz Odrzygóźdź, Piotr Miłoś, Yuhuai Wu, and Mateja Jamnik. Thor: Wielding hammers to integrate language models and automated theorem provers. *Advances in Neural Information Processing Systems*, 35:8360–8373, Nov 2022.
- [342] Fengqing Jiang. Identifying and mitigating vulnerabilities in llm-integrated applications. Master’s thesis, University of Washington, Jul 2024.
- [343] Nan Jiang, Shanchao Liang, Chengxiao Wang, Jiannan Wang, and Lin Tan. Latte: Improving latex recognition for tables and formulae with iterative refinement. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 4030–4038, Apr 2025.
- [344] Xiaowei Jiang, Liang Ou, Yanan Chen, Na Ao, Yu-Cheng Chang, Thomas Do, and Chin-Teng Lin. A fuzzy logic-based approach to predict human interaction by functional near-infrared spectroscopy. *IEEE Transactions on Fuzzy Systems*, Jan 2025.
- [345] Xue Jiang, Dezhen Xue, William Yi Wang, Jianjun Liu, Mingli Yang, Yanjing Su, et al. Ai4materials: Transforming the landscape of materials science and engineering. *Review of Materials Research*, page 100010, Jan 2025.
- [346] Licheng Jiao, Xue Song, Chao You, Xu Liu, Lingling Li, Puhua Chen, Xu Tang, Zhixi Feng, Fang Liu, Yuwei Guo, et al. Ai meets physics: a comprehensive survey. *Artificial Intelligence Review*, 57(9):256, Aug 2024.
- [347] Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. SWE-bench: Can language models resolve real-world github issues? In *The Twelfth International Conference on Learning Representations*, Jan 2024. URL <https://openreview.net/forum?id=VTF8yNQm66>.
- [348] Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William W Cohen, and Xinghua Lu. Pubmedqa: A dataset for biomedical research question answering. *arXiv preprint arXiv:1909.06146*, 2019.

- [349] Yiqiao Jin, Qinlin Zhao, Yiyang Wang, Hao Chen, Kaijie Zhu, Yijia Xiao, and Jindong Wang. Agentreview: Exploring peer review dynamics with llm agents. *arXiv preprint arXiv:2406.12708*, 2024.
- [350] Liqiang Jing, Zhehui Huang, Xiaoyang Wang, Wenlin Yao, Wenhao Yu, Kaixin Ma, Hongming Zhang, Xinya Du, and Dong Yu. Dsbench: How far are data science agents to becoming data science experts? *arXiv preprint arXiv:2409.07703*, 2024.
- [351] Wayne Johnson and Devon Proudfoot. Greater variability in judgements of the value of novel ideas. *Nature Human Behaviour*, 8(3):471–479, Jan 2024.
- [352] Cameron R Jones and Benjamin K Bergen. Large language models pass the turing test. *arXiv preprint arXiv:2503.23674*, 2025.
- [353] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Židek, Anna Potapenko, et al. Highly accurate protein structure prediction with alphafold. *nature*, 596(7873):583–589, Jul 2021.
- [354] Prerna Juneja and Tanushree Mitra. Human and technological infrastructures of fact-checking. *Proceedings of the ACM on Human-Computer Interaction*, 6(CSCW2):1–36, Nov 2022.
- [355] Rebecca E Kaiser, Farshad Sadr, Trevor Yuen, Till Krenz, Lee Chin-Chin, C Dominguez Sheela, Daru LL Ransford, and Erin Kobetz. 376 using a large language model to create lay summaries of clinical study descriptions. *Journal of Clinical and Translational Science*, 9(s1):116–116, Apr 2025.
- [356] Dmitry Kalashnikov, Alex Irpan, Peter Pastor, Julian Ibarz, Alexander Herzog, Eric Jang, Deirdre Quillen, Ethan Holly, Mrinal Kalakrishnan, Vincent Vanhoucke, et al. Scalable deep reinforcement learning for vision-based robotic manipulation. In *Conference on robot learning*, pages 651–673. PMLR, Oct 2018.
- [357] Subbarao Kambhampati, Karthik Valmeekam, Lin Guan, Mudit Verma, Kaya Stechly, Siddhant Bhambri, Lucas Paul Saldyt, and Anil B Murthy. Position: Llms can’t plan, but can help planning in llm-modulo frameworks. In *Forty-first International Conference on Machine Learning*, May 2024.
- [358] Dongyeop Kang, Waleed Ammar, Bhavana Dalvi, Madeleine van Zuylen, Sebastian Kohlmeier, Eduard Hovy, and Roy Schwartz. A dataset of peer reviews (peerread): Collection, insights and nlp applications. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1647–1661, Apr 2018.
- [359] Hyeonsu B Kang, Rafal Kocielnik, Andrew Head, Jiangjiang Yang, Matt Latzke, Aniket Kittur, Daniel S Weld, Doug Downey, and Jonathan Bragg. From who you know to what you read: Augmenting scientific recommendations with implicit social networks. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, pages 1–23, Apr 2022.
- [360] Hyeonsu B Kang, Nouran Soliman, Matt Latzke, Joseph Chee Chang, and Jonathan Bragg. Comlittee: Literature discovery with personal elected author committees. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, pages 1–20, Apr 2023.
- [361] Sungmin Kang, Bei Chen, Shin Yoo, and Jian-Guang Lou. Explainable automated debugging via large language model-driven scientific debugging. *Empirical Software Engineering*, 30(2):1–28, Mar 2025.

- [362] Chia Hsiang Kao, Wenting Zhao, Shreelekha Revankar, Samuel Speas, Snehal Bhagat, Rajeev Datta, Cheng Perng Phoo, Utkarsh Mall, Carl Vondrick, Kavita Bala, et al. Towards llm agents for earth observation. *arXiv preprint arXiv:2504.12110*, 2025.
- [363] Wei-Yu Kao and An-Zi Yen. Magic: Multi-argument generation with self-refinement for domain generalization in automatic fact-checking. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 10891–10902, May 2024.
- [364] Maryna Kapitonova and Tonio Ball. Human-ai teaming using large language models: Boosting brain-computer interfacing (bci) and brain research. *arXiv preprint arXiv:2501.01451*, 2024.
- [365] Andres Karjus. Machine-assisted quantizing designs: augmenting humanities and social sciences with artificial intelligence. *Humanities and Social Sciences Communications*, 12(1):1–18, Feb 2025.
- [366] Tetsu Kasanishi, Masaru Isonuma, Junichiro Mori, and Ichiro Sakata. Scireviewgen: a large-scale dataset for automatic literature review generation. *arXiv preprint arXiv:2305.15186*, 2023.
- [367] Uri Katz, Mosh Levy, and Yoav Goldberg. Knowledge navigator: Llm-guided browsing framework for exploratory search in scientific literature. *arXiv preprint arXiv:2408.15836*, 2024.
- [368] Nan Rosemary Ke, Danny P Sawyer, Hubert Soyer, Martin Engelcke, David P Reichert, Drew A Hudson, John Reid, Alexander Lerchner, Danilo Jimenez Rezende, Timothy P Lillicrap, et al. Can foundation models actively gather information in interactive environments to test hypotheses? *arXiv preprint arXiv:2412.06438*, 2024.
- [369] Martin Kerwer, Anita Chasiotis, Johannes Stricker, Armin Günther, and Tom Rosman. Straight from the scientist’s mouth—plain language summaries promote laypeople’s comprehension and knowledge acquisition when reading about individual research findings in psychology. *Collabra: Psychology*, 7(1), Feb 2021.
- [370] Farhana Keya, Gollam Rabby, Prasenjit Mitra, Sahar Vahdati, Sören Auer, and Yaser Jaradeh. Sci-idea: Context-aware scientific ideation using token and sentence embeddings. *arXiv preprint arXiv:2503.19257*, 2025.
- [371] Mohamed Khalifa and Mona Albadawy. Using artificial intelligence in academic writing and research: An essential productivity tool. *Computer Methods and Programs in Biomedicine Update*, page 100145, Mar 2024.
- [372] Mohamed Khalifa, Mona Albadawy, and Usman Iqbal. Advancing clinical decision support: The role of artificial intelligence across six domains. *Computer Methods and Programs in Biomedicine Update*, 5: 100142, Feb 2024.
- [373] Wasif Khan, Seowung Leem, Kyle B See, Joshua K Wong, Shaoting Zhang, and Ruogu Fang. A comprehensive survey of foundation models in medicine. *IEEE Reviews in Biomedical Engineering*, Jan 2025.
- [374] Jaeho Kim, Yunseok Lee, and Seulki Lee. Position: The ai conference peer review crisis demands author feedback and reviewer rewards. *arXiv preprint arXiv:2505.04966*, 2025.

- [375] Jaeyoung Kim, Jongho Lee, Hong-Jun Choi, Ting-Yao Hsu, Chieh-Yang Huang, Sungchul Kim, Ryan Rossi, Tong Yu, Clyde Lee Giles, Ting-Hao'Kenneth' Huang, et al. Multi-llm collaborative caption generation in scientific documents. *arXiv preprint arXiv:2501.02552*, 2025.
- [376] Jiho Kim, Sungjin Park, Yeonsu Kwon, Yohan Jo, James Thorne, and Edward Choi. Factkg: Fact verification via reasoning on knowledge graphs. *arXiv preprint arXiv:2305.06590*, 2023.
- [377] Min-Woo Kim, Hyo-Bin Park, Hee-Jin Ahn, Woo-Ram Park, Jae-Wan Jeon, Kyong-Ha Lee, Ryong Lee, and Dong-Geol Choi. Autopaperbench: An mllm-based framework for automatic generation of paper understanding evaluation benchmarks. *Electronics*, 14(6):1175, Mar 2025.
- [378] Seonok Kim. Medbiolm: Optimizing medical and biological qa with fine-tuned large language models and retrieval-augmented generation. *arXiv preprint arXiv:2502.03004*, 2025.
- [379] Ross D King, Kenneth E Whelan, Ffion M Jones, Philip GK Reiser, Christopher H Bryant, Stephen H Muggleton, Douglas B Kell, and Stephen G Oliver. Functional genomic hypothesis generation and experimentation by a robot scientist. *Nature*, 427(6971):247–252, Jan 2004.
- [380] Chhavi Kirtani, Madhav Krishan Garg, Tejash Prasad, Tanmay Singhal, Murari Mandal, and Dhruv Kumar. Revieweval: An evaluation framework for ai-generated reviews. *arXiv preprint arXiv:2502.11736*, 2025.
- [381] Stephen T Knox, Kai E Wu, Nazrul Islam, Roisin O'Connell, Peter M Pittaway, Kudakwashe E Chingono, John Oyekan, George Panoutsos, Thomas W Chamberlain, Richard A Bourne, et al. Self-driving laboratory platform for many-objective self-optimisation of polymer nanoparticle synthesis with cloud-integrated machine learning and orthogonal online analytics. *Polymer Chemistry*, 16(12):1355–1364, Feb 2025.
- [382] Patrick Tser Jern Kon, Jiachen Liu, Qiuyi Ding, Yiming Qiu, Zhenning Yang, Yibo Huang, Jayanth Srinivasa, Myungjin Lee, Mosharaf Chowdhury, and Ang Chen. Curie: Toward rigorous and automated scientific experimentation with ai agents. *arXiv preprint arXiv:2502.16069*, 2025.
- [383] Patrick Tser Jern Kon, Jiachen Liu, Xinyi Zhu, Qiuyi Ding, Jingjia Peng, Jiarong Xing, Yibo Huang, Yiming Qiu, Jayanth Srinivasa, Myungjin Lee, et al. Exp-bench: Can ai conduct ai research experiments? *arXiv preprint arXiv:2505.24785*, 2025.
- [384] Kayvan Kousha and Mike Thelwall. Artificial intelligence to support publishing and peer review: A summary and review. *Learned Publishing*, 37(1):4–12, Aug 2024.
- [385] Mario Krenn, Lorenzo Buffoni, Bruno Coutinho, Sagi Eppel, Jacob Gates Foster, Andrew Gritsevskiy, Harlin Lee, Yichao Lu, Joao P Moutinho, Nima Sanjabi, et al. Predicting the future of ai with ai: High-quality link prediction in an exponentially growing knowledge network. *arXiv preprint arXiv:2210.00881*, 2022.
- [386] Christin Katharina Kreutz and Ralf Schenkel. Scientific paper recommendation systems: a literature review of recent publications. *International journal on digital libraries*, 23(4):335–369, Oct 2022.
- [387] Amrith Krishna, Sebastian Riedel, and Andreas Vlachos. Proofver: Natural logic theorem proving for fact verification. *Transactions of the Association for Computational Linguistics*, 10:1013–1030, Sep 2022.

- [388] Agustinus Kristiadi, Felix Strieth-Kalthoff, Marta Skreta, Pascal Poupart, Alán Aspuru-Guzik, and Geoff Pleiss. A sober look at llms for material discovery: Are they actually good for bayesian optimization over molecules? *arXiv preprint arXiv:2402.05015*, 2024.
- [389] Anastasia Krithara, Anastasios Nentidis, Konstantinos Bougiatiotis, and Georgios Paliouras. Bioasq-qa: A manually curated corpus for biomedical question answering. *Scientific Data*, 10(1):170, Mar 2023.
- [390] Max Ku, Thomas Chong, Jonathan Leung, Krish Shah, Alvin Yu, and Wenhui Chen. Theoremexplainer: Towards video-based multimodal explanations for llm theorem understanding. *arXiv preprint arXiv:2502.19400*, 2025.
- [391] Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. Clam: Selective clarification for ambiguous questions with generative language models. *arXiv preprint arXiv:2212.07769*, 2022.
- [392] Adithya Kulkarni, Fatimah Alotaibi, Xinyue Zeng, Longfeng Wu, Tong Zeng, Barry Menglong Yao, Minqian Liu, Shuaicheng Zhang, Lifu Huang, and Dawei Zhou. Scientific hypothesis generation and validation: Methods, datasets, and future directions. *arXiv preprint arXiv:2505.04651*, 2025.
- [393] Asheesh Kumar, Tirthankar Ghosal, Saprativa Bhattacharjee, and Asif Ekbal. Towards automated meta-review generation via an nlp/ml pipeline in different stages of the scholarly peer review process. *International Journal on Digital Libraries*, 25(3):493–504, Apr 2024.
- [394] Harsh Kumar, Jonathan Vincentius, Ewan Jordan, and Ashton Anderson. Human creativity in the age of llms: Randomized experiments on divergent and convergent thinking. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*, pages 1–18, Apr 2025.
- [395] Sandeep Kumar, Tirthankar Ghosal, and Asif Ekbal. When reviewers lock horn: Finding disagreement in scientific peer reviews. *arXiv preprint arXiv:2310.18685*, 2023.
- [396] Sandeep Kumar, Tirthankar Ghosal, Vinayak Goyal, and Asif Ekbal. Can large language models unlock novel scientific research ideas? *arXiv preprint arXiv:2409.06185*, 2024.
- [397] Sujit Kumar, Anshul Sharma, Siddharth Hemant Khincha, Gargi Shroff, Sanasam Ranbir Singh, and Rahul Mishra. Sciclamhunt: A large dataset for evidence-based scientific claim verification. *arXiv preprint arXiv:2502.10003*, 2025.
- [398] Tsung-Ting Kuo, Rodney A Gabriel, Jejo Koola, Robert T Schooley, and Lucila Ohno-Machado. Distributed cross-learning for equitable federated models-privacy-preserving prediction on data from five california hospitals. *Nature Communications*, 16(1):1371, Feb 2025.
- [399] Ilia Kuznetsov, Osama Mohammed Afzal, Koen Dercksen, Nils Dycke, Alexander Goldberg, Tom Hope, Dirk Hovy, Jonathan K Kummerfeld, Anne Lauscher, Kevin Leyton-Brown, et al. What can natural language processing do for peer review? *arXiv preprint arXiv:2405.06563*, 2024.
- [400] J Kvpil, G Borca-Tasciuc, H Bossi, K Chen, Y Chen, Y Corrales Morales, H Da Costa, C Da Silva, C Dean, J Durham, et al. Intelligent experiments through real-time ai: Fast data processing and autonomous detector control for sphenix and future eic detectors. *arXiv preprint arXiv:2501.04845*, 2025.
- [401] Daeun Kyung, Hyunseung Chung, Seongsu Bae, Jiho Kim, Jae Ho Sohn, Taerim Kim, Soo Kyung Kim, and Edward Choi. Patientsim: A persona-driven simulator for realistic doctor-patient interactions. *arXiv preprint arXiv:2505.17818*, 2025.

- [402] Arash Lagzian, Srinivas Anumasa, and Dianbo Liu. Multi-novelty: Improve the diversity and novelty of contents generated by large language models via inference-time multi-views brainstorming. *arXiv preprint arXiv:2502.12700*, 2025.
- [403] Yuxuan Lai, Yupeng Wu, Yidan Wang, Wenpeng Hu, and Chen Zheng. Instruct large language models to generate scientific literature survey step by step. In *CCF International Conference on Natural Language Processing and Chinese Computing*, pages 484–496. Springer, Nov 2024.
- [404] Zheyuan Lai and Yingming Pu. Prim: Principle-inspired material discovery through multi-agent collaboration. *arXiv preprint arXiv:2504.08810*, 2025.
- [405] Jakub Lála, Odhran O’Donoghue, Aleksandar Shtedritski, Sam Cox, Samuel G Rodriques, and Andrew D White. Paperqa: Retrieval-augmented generative agent for scientific research. *arXiv preprint arXiv:2312.07559*, 2023.
- [406] Hilbert Yuen In Lam, Xing Er Ong, and Marek Mutwil. Large language models in plant biology. *Trends in Plant Science*, Oct 2024.
- [407] Guillaume Lample, Timothee Lacroix, Marie-Anne Lachaux, Aurelien Rodriguez, Amaury Hayat, Thibaut Lavril, Gabriel Ebner, and Xavier Martinet. Hypertree proof search for neural theorem proving. *Advances in Neural Information Processing Systems*, 35:26337–26349, Nov 2022.
- [408] Ge Lan, Mengting Hu, Ye Li, and Yuzhi Zhang. Contrastive knowledge integrated graph neural networks for chinese medical text classification. *Engineering Applications of Artificial Intelligence*, 122: 106057, Jun 2023.
- [409] Wuyang Lan, Wenzheng Wang, Changwei Ji, Guoxing Yang, Yongbo Zhang, Xiaohong Liu, Song Wu, and Guangyu Wang. Clinicalgpt-r1: Pushing reasoning capability of generalist disease diagnosis with large language model. *arXiv preprint arXiv:2504.09421*, 2025.
- [410] Zhiqian Lan, Yuxuan Jiang, Ruiqi Wang, Xuanbing Xie, Rongkui Zhang, Yicheng Zhu, Peihang Li, Tianshuo Yang, Tianxing Chen, Haoyu Gao, et al. Autobio: A simulation and benchmark for robotic automation in digital biology laboratory. *arXiv preprint arXiv:2505.14030*, 2025.
- [411] Robert Tjarko Lange, Aaditya Prasad, Qi Sun, Maxence Faldor, Yujin Tang, and David Ha. The ai cuda engineer: Agentic cuda kernel discovery, optimization and composition. Technical report, Technical report, Sakana AI, 02 2025, Feb 2025.
- [412] Antonio Laverghetta Jr, Tuhin Chakrabarty, Tom Hope, Jimmy Pronchick, Krupa Bhawsar, and Roger E Beaty. How do humans and language models reason about creativity? a comparative analysis. *arXiv preprint arXiv:2502.03253*, 2025.
- [413] Nam Le Hai, Dung Manh Nguyen, and Nghi DQ Bui. Repoexec: Evaluate code generation with a repository-level executable benchmark. *arXiv e-prints*, pages arXiv–2406, 2024.
- [414] Jingoo Lee, Kyungho Lim, Young-Chul Jung, and Byung-Hoon Kim. Psyche: A multi-faceted patient simulation framework for evaluation of psychiatric assessment conversational agents. *arXiv preprint arXiv:2501.01594*, 2025.
- [415] Jisoo Lee, Jieun Lee, and Jeong-Ju Yoo. The role of large language models in the peer-review process: opportunities and challenges for medical journal reviewers and editors. *Journal of Educational Evaluation for Health Professions*, 22, Jan 2025.

- [416] Mina Lee, Percy Liang, and Qian Yang. Coauthor: Designing a human-ai collaborative writing dataset for exploring language model capabilities. In *Proceedings of the 2022 CHI conference on human factors in computing systems*, pages 1–19, Apr 2022.
- [417] Namkyeong Lee, Edward De Brouwer, Ehsan Hajiramezanali, Tommaso Biancalani, Chanyoung Park, and Gabriele Scalia. Rag-enhanced collaborative llm agents for drug discovery. *arXiv preprint arXiv:2502.17506*, 2025.
- [418] Seungyeon Lee, Ruoqi Liu, Feixiong Cheng, and Ping Zhang. A deep subgrouping framework for precision drug repurposing via emulating clinical trials on real-world patient data. *arXiv preprint arXiv:2412.20373*, 2024.
- [419] Suk Ki Lee and Hyunwoong Ko. Generative machine learning in adaptive control of dynamic manufacturing processes: A review. *arXiv preprint arXiv:2505.00210*, 2025.
- [420] Yoonjoo Lee, Hyeonsu B Kang, Matt Latzke, Juho Kim, Jonathan Bragg, Joseph Chee Chang, and Pao Siangliulue. Paperweaver: Enriching topical paper alerts by contextualizing recommended papers with user-collected papers. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, pages 1–19, May 2024.
- [421] Steven A Lehr, Aylin Caliskan, Suneragiri Liyanage, and Mahzarin R Banaji. Chatgpt as research scientist: probing gpt’s capabilities as a research librarian, research ethicist, data generator, and data predictor. *Proceedings of the National Academy of Sciences*, 121(35):e2404328121, Jul 2024.
- [422] Yan Leng, Yi Zhong, Zhi Gu, Peiyi Li, Haoting Cui, Xing Li, Yang Liu, and Jiayu Wan. Intelligent, personalized scientific assistant via large language models for solid-state battery research. *ACS Materials Letters*, 7(5):1807–1816, Apr 2025.
- [423] Shi Xuan Leong, Sergio Pablo-García, Brandon Wong, and Alán Aspuru-Guzik. Mermaid: Universal multimodal mining of chemical reactions from pdfs using vision-language models. Mar 2025.
- [424] Sergey Levine, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. End-to-end training of deep visuomotor policies. *Journal of Machine Learning Research*, 17(39):1–40, Apr 2016.
- [425] Sergey Levine, Peter Pastor, Alex Krizhevsky, Julian Ibarz, and Deirdre Quillen. Learning hand-eye coordination for robotic grasping with deep learning and large-scale data collection. *The International journal of robotics research*, 37(4-5):421–436, Jun 2018.
- [426] Kevin Leyton-Brown, Yatin Nandwani, Hedayat Zarkoob, Chris Cameron, Neil Newman, Dinesh Raghu, et al. Matching papers and reviewers at large conferences. *Artificial Intelligence*, 331:104119, Jun 2024.
- [427] Bohan Li, Jiannan Guan, Longxu Dou, Yunlong Feng, Dingzirui Wang, Yang Xu, Enbo Wang, Qiguang Chen, Bichen Wang, Xiao Xu, et al. Can large language models understand you better? an mbti personality detection dataset aligned with population traits. *arXiv preprint arXiv:2412.12510*, 2024.
- [428] Chuhan Li, Ziyao Shangguan, Yilun Zhao, Deyuan Li, Yixin Liu, and Arman Cohan. M3SciQA: A multi-modal multi-document scientific QA benchmark for evaluating foundation models. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 15419–15446, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-emnlp.904. URL <https://aclanthology.org/2024.findings-emnlp.904/>.

- [429] Junkai Li, Yungchwei Lai, Weitao Li, Jingyi Ren, Meng Zhang, Xinhui Kang, Siyu Wang, Peng Li, Ya-Qin Zhang, Weizhi Ma, et al. Agent hospital: A simulacrum of hospital with evolvable medical agents. *arXiv preprint arXiv:2405.02957*, 2024.
- [430] Kun Li, Zhennan Wu, Shoupeng Wang, and Wenbin Hu. Drugpilot: Llm-based parameterized reasoning agent for drug discovery. *arXiv preprint arXiv:2505.13940*, 2025.
- [431] Lei Li, Yuqi Wang, Runxin Xu, Peiyi Wang, Xiachong Feng, Lingpeng Kong, and Qi Liu. Multimodal ArXiv: A dataset for improving scientific comprehension of large vision-language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14369–14387, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.775. URL <https://aclanthology.org/2024.acl-long.775/>.
- [432] Li Li, Yuxi Fan, Mike Tse, and Kuo-Yi Lin. A review of applications in federated learning. *Computers & Industrial Engineering*, 149:106854, Nov 2020.
- [433] Long Li, Weiwen Xu, Jiayan Guo, Ruochen Zhao, Xingxuan Li, Yuqian Yuan, Boqiang Zhang, Yuming Jiang, Yifei Xin, Ronghao Dang, et al. Chain of ideas: Revolutionizing research via novel idea development with llm agents. *arXiv preprint arXiv:2410.13185*, 2024.
- [434] Miao Li, Eduard Hovy, and Jey Han Lau. Summarizing multiple documents with conversational structure for meta-review generation. *arXiv preprint arXiv:2305.01498*, 2023.
- [435] Miao Li, Jey Han Lau, and Eduard Hovy. A sentiment consolidation framework for meta-review generation. *arXiv preprint arXiv:2402.18005*, 2024.
- [436] Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, et al. Starcoder: may the source be with you! *arXiv preprint arXiv:2305.06161*, 2023.
- [437] Ruifeng Li, Dongzhan Zhou, Ancheng Shen, Ao Zhang, Mao Su, Mingqian Li, Hongyang Chen, Gang Chen, Yin Zhang, Shufei Zhang, et al. Physical formula enhanced multi-task learning for pharmacokinetics prediction. *arXiv preprint arXiv:2404.10354*, 2024.
- [438] Ruifeng Li, Mingqian Li, Wei Liu, Yuhua Zhou, Xiangxin Zhou, Yuan Yao, Qiang Zhang, and Hongyang Chen. Unimatch: Universal matching from atom to task for few-shot drug discovery. *arXiv preprint arXiv:2502.12453*, 2025.
- [439] Ruikun Li, Yan Lu, Shixiang Tang, Biqing Qi, and Wanli Ouyang. Mllm-based discovery of intrinsic coordinates and governing equations from high-dimensional data. *arXiv preprint arXiv:2505.11940*, 2025.
- [440] Ruochen Li, Liqiang Jing, Chi Han, Jiawei Zhou, and Xinya Du. Learning to generate research idea with dynamic control. *arXiv preprint arXiv:2412.14626*, 2024.
- [441] Ruochen Li, Teerth Patel, Qingyun Wang, and Xinya Du. Mlr-copilot: Autonomous machine learning research based on large language models agents. *arXiv preprint arXiv:2408.14033*, 2024.
- [442] Ruosen Li, Ruochen Li, Barry Wang, and Xinya Du. Iqa-eval: Automatic evaluation of human-model interactive question answering. *Advances in Neural Information Processing Systems*, 37:109894–109921, Dec 2024.

- [443] Sitong Li, Stefano Padilla, Pierre Le Bras, Junyu Dong, and Mike Chantler. A review of llm-assisted ideation. *arXiv preprint arXiv:2503.00946*, 2025.
- [444] Wenyu Li, Zhitao Mao, Zhengyang Xiao, Xiaoping Liao, Mattheos Koffas, Yixin Chen, Hongwu Ma, and Yinjie J Tang. Large language model for knowledge synthesis and ai-enhanced biomanufacturing. *Trends in Biotechnology*, Mar 2025.
- [445] Xiangci Li and Jessica Ouyang. Related work and citation text generation: A survey. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 13846–13864, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.767. URL <https://aclanthology.org/2024.emnlp-main.767/>.
- [446] Xiangci Li and Jessica Ouyang. Explaining relationships among research papers. *arXiv preprint arXiv:2402.13426*, 2024.
- [447] Xiangci Li and Jessica Ouyang. Related work and citation text generation: A survey. *arXiv preprint arXiv:2404.11588*, 2024.
- [448] Xiangci Li, Yi-Hui Lee, and Jessica Ouyang. Cited text spans for citation text generation. *arXiv preprint arXiv:2309.06365*, 2023.
- [449] Xiangyu Li and Jingqiang Chen. Scirgc: Multi-granularity citation recommendation and citation sentence preference alignment. *arXiv preprint arXiv:2505.20103*, 2025.
- [450] Xiaobo Li, Yu Che, Linjiang Chen, Tao Liu, Kewei Wang, Lunjie Liu, Haofan Yang, Edward O Pyzer-Knapp, and Andrew I Cooper. Sequential closed-loop bayesian optimization as a guide for organic molecular metallophotocatalyst formulation discovery. *Nature Chemistry*, 16(8):1286–1294, Jun 2024.
- [451] Yifei Li, Hanane Nour Moussa, Ziru Chen, Shijie Chen, Botao Yu, Mingyi Xue, Benjamin Burns, Tzu-Yao Chiu, Vishal Dey, Zitong Lu, et al. Autosdt: Scaling data-driven discovery tasks toward open co-scientists. *arXiv preprint arXiv:2506.08140*, 2025.
- [452] Yugang Li, Baizhou Wu, Yuqi Huang, and Shenghua Luan. Developing trustworthy artificial intelligence: insights from research on interpersonal, human-automation, and human-ai trust. *Frontiers in Psychology*, 15:1382693, Apr 2024.
- [453] Zekun Li, Xianjun Yang, Kyuri Choi, Wanrong Zhu, Ryan Hsieh, HyeonJung Kim, Jin Hyuk Lim, Sungyoung Ji, Byungju Lee, Xifeng Yan, et al. Mmsci: A dataset for graduate-level multi-discipline multimodal scientific understanding. *arXiv preprint arXiv:2407.04903*, 2024.
- [454] Zhaoyu Li, Jialiang Sun, Logan Murphy, Qidong Su, Zenan Li, Xian Zhang, Kaiyu Yang, and Xujie Si. A survey on deep learning for theorem proving. *arXiv preprint arXiv:2404.09939*, 2024.
- [455] Zhi Li and Xiaozhu Zou. A review on personalized academic paper recommendation. *Comput. Inf. Sci.*, 12(1):33–43, Jan 2019.
- [456] Zhuofan Li and Corey M Abramson. Ethnography and machine learning: Synergies and new directions. *arXiv preprint arXiv:2412.06087*, 2024.

- [457] Ziming Li, Qianbo Zang, David Ma, Jiawei Guo, Tuney Zheng, Minghao Liu, Xinyao Niu, Yue Wang, Jian Yang, Jiaheng Liu, et al. Autokaggle: A multi-agent framework for autonomous data science competitions. *arXiv preprint arXiv:2410.20424*, 2024.
- [458] Haotong Liang, Chuangye Wang, Heshan Yu, Dylan Kirsch, Rohit Pant, Austin McDannald, A Gilad Kusne, Ji-Cheng Zhao, and Ichiro Takeuchi. Real-time experiment-theory closed-loop interaction for autonomous materials science. *arXiv preprint arXiv:2410.17430*, 2024.
- [459] Jing Liang. The application of artificial intelligence-assisted technology in cultural and creative product design. *Scientific Reports*, 14(1):31069, Dec 2024.
- [460] Siting Liang and Daniel Sonntag. Explainable biomedical claim verification with large language models. *arXiv preprint arXiv:2502.21014*, 2025.
- [461] Weixin Liang, Yuhui Zhang, Hancheng Cao, Binglu Wang, Daisy Yi Ding, Xinyu Yang, Kailas Vodrahalli, Siyu He, Daniel Scott Smith, Yian Yin, et al. Can large language models provide useful feedback on research papers? a large-scale empirical analysis. *NEJM AI*, 1(8):AIoa2400196, Jul 2024.
- [462] Xun Liang, Jiawei Yang, Yezhaohui Wang, Chen Tang, Zifan Zheng, Shichao Song, Zehao Lin, Yebin Yang, Simin Niu, Hanyu Wang, et al. Surveyx: Academic survey automation via large language models. *arXiv preprint arXiv:2502.14776*, 2025.
- [463] Zhenwen Liang, Kehan Guo, Gang Liu, Taicheng Guo, Yujun Zhou, Tianyu Yang, Jiajun Jiao, Renjie Pi, Jipeng Zhang, and Xiangliang Zhang. SceMQA: A scientific college entrance level multimodal question answering benchmark. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 109–119, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-short.11. URL <https://aclanthology.org/2024.acl-short.11/>.
- [464] Zhehui Liao, Maria Antoniak, Inyoung Cheong, Evie Yu-Yen Cheng, Ai-Heng Lee, Kyle Lo, Joseph Chee Chang, and Amy X Zhang. Llm as research tools: A large scale survey of researchers’ usage and perceptions. *arXiv preprint arXiv:2411.05025*, 2024.
- [465] Daniel J. Liebling, Malcolm Kane, Madeleine Grunde-McLaughlin, Ian Lang, Subhashini Venugopalan, and Michael Brenner. Towards AI-assisted academic writing. In Peter Jansen, Bhavana Dalvi Mishra, Harsh Trivedi, Bodhisattwa Prasad Majumder, Tom Hope, Tushar Khot, Doug Downey, and Eric Horvitz, editors, *Proceedings of the 1st Workshop on AI and Scientific Discovery: Directions and Opportunities*, pages 31–45, Albuquerque, New Mexico, USA, May 2025. Association for Computational Linguistics. ISBN 979-8-89176-224-4. URL <https://aclanthology.org/2025.aisd-main.4/>.
- [466] Cong William Lin and Wu Zhu. Divergent llm adoption and heterogeneous convergence paths in research writing. *arXiv preprint arXiv:2504.13629*, 2025.
- [467] Ethan Lin, Zhiyuan Peng, and Yi Fang. Evaluating and enhancing large language models for novelty assessment in scholarly publications. In Peter Jansen, Bhavana Dalvi Mishra, Harsh Trivedi, Bodhisattwa Prasad Majumder, Tom Hope, Tushar Khot, Doug Downey, and Eric Horvitz, editors, *Proceedings of the 1st Workshop on AI and Scientific Discovery: Directions and Opportunities*, pages 46–57, Albuquerque, New Mexico, USA, May 2025. Association for Computational Linguistics. ISBN 979-8-89176-224-4. URL <https://aclanthology.org/2025.aisd-main.5/>.

- [468] Haohan Lin, Zhiqing Sun, Sean Welleck, and Yiming Yang. Lean-star: Learning to interleave thinking and proving. *arXiv preprint arXiv:2407.10040*, 2024.
- [469] Jialiang Lin, Jiaxin Song, Zhangping Zhou, Yidong Chen, and Xiaodong Shi. Automated scholarly paper review: Concepts, technologies, and challenges. *Information fusion*, 98:101830, Oct 2023.
- [470] Jialiang Lin, Jiaxin Song, Zhangping Zhou, Yidong Chen, and Xiaodong Shi. Moprdr: A multidisciplinary open peer review dataset. *Neural Computing and Applications*, 35(34):24191–24206, Sep 2023.
- [471] Tung-Wei Lin, Runing Yang, Zain ul Abdeen, Alberto Sangiovanni-Vincentelli, Haibo Huang, and Ming Jin. Llm tackle meta-analysis: Automating scientific hypothesis generation with statistical rigor. In *2nd AI4Research Workshop: Towards a Knowledge-grounded Scientific Research Lifecycle*, Dec 2024.
- [472] Xinna Lin, Siqi Ma, Junjie Shan, Xiaojing Zhang, Shell Xu Hu, Tiannan Guo, Stan Z Li, and Kaicheng Yu. Biokgbench: A knowledge graph checking benchmark of ai agent for biomedical science. *arXiv preprint arXiv:2407.00466*, 2024.
- [473] Xuan Lin, Qingrui Liu, Hongxin Xiang, Daojian Zeng, and Xiangxiang Zeng. Enhancing chemical reaction and retrosynthesis prediction with large language model and dual-task learning. *arXiv preprint arXiv:2505.02639*, 2025.
- [474] Xule Lin. Cognition emerges: Agency, dimensions, and dynamics in human-ai knowledge co-creation. *arXiv preprint arXiv:2505.03105*, 2025.
- [475] Zhicheng Lin. Beyond principlism: practical strategies for ethical ai use in research practices. *AI and Ethics*, pages 1–13, Oct 2024.
- [476] Zhicheng Lin. Techniques for supercharging academic writing with generative ai. *Nature Biomedical Engineering*, pages 1–6, Mar 2024.
- [477] Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024.
- [478] Bang Liu, Xinfeng Li, Jiayi Zhang, Jinlin Wang, Tanjin He, Sirui Hong, Hongzhang Liu, Shaokun Zhang, Kaitao Song, Kunlun Zhu, et al. Advances and challenges in foundation agents: From brain-inspired intelligence to evolutionary, collaborative, and safe systems. *arXiv preprint arXiv:2504.01990*, 2025.
- [479] Chengwei Liu, Chong Wang, Jiayue Cao, Jingquan Ge, Kun Wang, Lvy Zhang, Ming-Ming Cheng, Penghai Zhao, Tianlin Li, Xiaojun Jia, et al. A vision for auto research with llm agents. *arXiv preprint arXiv:2504.18765*, 2025.
- [480] Chengwu Liu, Jianhao Shen, Huajian Xin, Zhengying Liu, Ye Yuan, Haiming Wang, Wei Ju, Chuanyang Zheng, Yichun Yin, Lin Li, et al. Fimo: A challenge formal dataset for automated theorem proving. *arXiv preprint arXiv:2309.04295*, 2023.
- [481] Chunwei Liu, Enrique Noriega-Atala, Adarsh Pyarelal, Clayton T Morrison, and Mike Cafarella. Variable extraction for model recovery in scientific literature. In Peter Jansen, Bhavana Dalvi Mishra, Harsh Trivedi, Bodhisattwa Prasad Majumder, Tom Hope, Tushar Khot, Doug Downey, and Eric Horvitz, editors, *Proceedings of the 1st Workshop on AI and Scientific Discovery: Directions and Opportunities*, pages 1–12, Albuquerque, New Mexico, USA, May 2025. Association for Computational Linguistics. ISBN 979-8-89176-224-4. doi: 10.18653/v1/2025.aisd-main.1. URL <https://aclanthology.org/2025.aisd-main.1/>.

- [482] Fengming Liu and Shubin Yu. Step further towards automated social science: An ai-powered interview platform. *Available at SSRN*, Apr 2024.
- [483] Haokun Liu, Yangqiaoyu Zhou, Mingxuan Li, Chenfei Yuan, and Chenhao Tan. Literature meets data: A synergistic approach to hypothesis generation. *arXiv preprint arXiv:2410.17309*, 2024.
- [484] Haokun Liu, Sicong Huang, Jingyu Hu, Yangqiaoyu Zhou, and Chenhao Tan. Hypobench: Towards systematic and principled benchmarking for hypothesis generation. *arXiv preprint arXiv:2504.11524*, 2025.
- [485] Haoyu Liu, Yifu Tang, Zizhao Zhang, Zeyu Zheng, and Tingyu Zhu. Large language model assisted experiment design with generative human-behavior agents. In *2024 Winter Simulation Conference (WSC)*, pages 2751–2762. IEEE, Dec 2024.
- [486] Jiachang Liu, Qi Zhang, Chongyang Shi, Usman Naseem, Shoujin Wang, and Ivor Tsang. Causal intervention for abstractive related work generation. *arXiv preprint arXiv:2305.13685*, 2023.
- [487] Jiaheng Liu, Dawei Zhu, Zhiqi Bai, Yancheng He, Huanxuan Liao, Haoran Que, Zekun Wang, Chenchen Zhang, Ge Zhang, Jiebin Zhang, et al. A comprehensive survey on long context language modeling. *arXiv preprint arXiv:2503.17407*, 2025.
- [488] Mengchen Liu, Qixiu Li, Dongdong Chen, Dong Chen, Jianmin Bao, and Yunsheng Li. Synchart: Synthesizing charts from language models. *arXiv preprint arXiv:2409.16517*, 2024.
- [489] Ruibo Liu, Ruixin Yang, Chenyan Jia, Ge Zhang, Denny Zhou, Andrew M Dai, Diyi Yang, and Soroush Vosoughi. Training socially aligned language models in simulated human society. *arXiv preprint arXiv:2305.16960*, 2, 2023.
- [490] Ryan Liu and Nihar B Shah. Reviewergpt? an exploratory study on using large language models for paper reviewing. *arXiv preprint arXiv:2306.00622*, 2023.
- [491] Shuaiqi Liu, Jiannong Cao, Ruosong Yang, and Zhiyuan Wen. Generating a structured summary of numerous academic papers: Dataset and method. *arXiv preprint arXiv:2302.04580*, 2023.
- [492] Siyi Liu, Chen Gao, and Yong Li. Large language model agent for hyper-parameter optimization. *arXiv preprint arXiv:2402.01881*, 2024.
- [493] Sizhe Liu, Yizhou Lu, Siyu Chen, Xiyang Hu, Jieyu Zhao, Yingzhou Lu, and Yue Zhao. Drugagent: Automating ai-aided drug discovery programming through llm multi-agent collaboration. *arXiv preprint arXiv:2411.15692*, 2024.
- [494] Wanhao Liu, Zonglin Yang, Jue Wang, Lidong Bing, Di Zhang, Dongzhan Zhou, Yuqiang Li, Houqiang Li, Erik Cambria, and Wanli Ouyang. Moose-chem3: Toward experiment-guided hypothesis ranking via simulated experimental feedback. *arXiv preprint arXiv:2505.17873*, 2025.
- [495] Xiao Liu, Xinyi Dong, Xinyang Gao, Yansong Feng, and Xun Pang. Improving research idea generation through data: An empirical investigation in social science. *arXiv preprint arXiv:2505.21396*, 2025.
- [496] Xiaochuan Liu, Ruihua Song, Xiting Wang, and Xu Chen. Select, read, and write: A multi-agent framework of full-text-based related work generation. *arXiv preprint arXiv:2505.19647*, 2025.
- [497] Yan Liu, Zonglin Yang, Soujanya Poria, Thanh-Son Nguyen, and Erik Cambria. Harnessing large language models for scientific novelty detection. *arXiv preprint arXiv:2505.24615*, 2025.

- [498] Yijun Liu, Jinzheng Yu, Yang Xu, Zhongyang Li, and Qingfu Zhu. A survey on transformer context extension: Approaches and evaluation. *arXiv preprint arXiv:2503.13299*, 2025.
- [499] Yiren Liu, Si Chen, Haocong Cheng, Mengxia Yu, Xiao Ran, Andrew Mo, Yiliu Tang, and Yun Huang. How ai processing delays foster creativity: Exploring research question co-creation with an llm-based agent. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, pages 1–25, May 2024.
- [500] Yujie Liu, Zonglin Yang, Tong Xie, Jinjie Ni, Ben Gao, Yuqiang Li, Shixiang Tang, Wanli Ouyang, Erik Cambria, and Dongzhan Zhou. Researchbench: Benchmarking llms in scientific discovery via inspiration-based task decomposition. *arXiv preprint arXiv:2503.21248*, 2025.
- [501] Yuyang Liu, Liuzhenghao Lv, Xiancheng Zhang, Li Yuan, and Yonghong Tian. Bioprobench: Comprehensive dataset and benchmark in biological protocol understanding and reasoning. *arXiv preprint arXiv:2505.07889*, 2025.
- [502] Zijun Liu, Kaiming Liu, Yiqi Zhu, Xuanyu Lei, Zonghan Yang, Zhenhe Zhang, Peng Li, and Yang Liu. Aigs: Generating science from ai-powered automated falsification. *arXiv preprint arXiv:2411.11910*, 2024.
- [503] Stanley Lo, Sterling G Baird, Joshua Schrier, Ben Blaiszik, Nessa Carson, Ian Foster, Andrés Aguilar-Granda, Sergei V Kalinin, Benji Maruyama, Maria Politi, et al. Review of low-cost self-driving laboratories in chemistry and materials science: the “frugal twin” concept. *Digital Discovery*, 3(5): 842–868, Feb 2024.
- [504] Joseph R Loffredo and Suyeol Yun. Agent-enhanced large language models for researching political institutions. *arXiv preprint arXiv:2503.13524*, 2025.
- [505] Renze Lou, Hanzi Xu, Sijia Wang, Jiangshu Du, Ryo Kamoi, Xiaoxin Lu, Jian Xie, Yuxuan Sun, Yusen Zhang, Jihyun Janice Ahn, et al. Aaar-1.0: Assessing ai’s potential to assist research. *arXiv preprint arXiv:2410.22394*, 2024.
- [506] Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, et al. Starcoder 2 and the stack v2: The next generation. *arXiv preprint arXiv:2402.19173*, 2024.
- [507] 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.
- [508] Cong Lu, Shengran Hu, and Jeff Clune. Beyond benchmarking: Automated capability discovery via model self-exploration. In *Language Gamification-NeurIPS 2024 Workshop*, Oct 2024.
- [509] Yingzhou Lu, Yaojun Hu, and Chenhao Li. Drugclip: Contrastive drug-disease interaction for drug repurposing. *arXiv preprint arXiv:2407.02265*, 2024.
- [510] Hewitt Luke, Ashokkumar Ashwini, Ghezae Isaias, and Willer Robb. Predicting results of social science experiments using large language models, Aug 2024. URL <https://samim.io/dl/Predicting%20results%20of%20social%20science%20experiments%20using%20large%20language%20models.pdf>. Predicting Results of Social Science Experiments Using Large Language Models.

- [511] Erpai Luo, Jinmeng Jia, Yifan Xiong, Xiangyu Li, Xiaobo Guo, Baoqi Yu, Lei Wei, and Xuegong Zhang. Benchmarking ai scientists in omics data-driven biological research. *arXiv preprint arXiv:2505.08341*, 2025.
- [512] Junyu Luo, Cheng Qian, Lucas Glass, and Fenglong Ma. Clinical trial retrieval via multi-grained similarity learning. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 2950–2954, Jul 2024.
- [513] Man Luo, Zikai Xie, Huirong Li, Baicheng Zhang, Jiaqi Cao, Yan Huang, Hang Qu, Qing Zhu, Linjiang Chen, Jun Jiang, et al. Physics-informed, dual-objective optimization of high-entropy-alloy nanozymes by a robotic ai chemist. *Matter*, 8(4), Apr 2025.
- [514] Shunyang Luo, Yuqi Tang, Mingyuan Jiang, Kehua Feng, Qiang Zhang, and Keyan Ding. Generating multiple choice questions from scientific literature via large language models. In *2024 IEEE International Conference on Knowledge Graph (ICKG)*, pages 219–226, Feb 2024. doi: 10.1109/ICKG63256.2024.00035.
- [515] Xiaoliang Luo, Akilles Rechartt, Guangzhi Sun, Kevin K Nejad, Felipe Yáñez, Bati Yilmaz, Kangjoo Lee, Alexandra O Cohen, Valentina Borghesani, Anton Pashkov, et al. Large language models surpass human experts in predicting neuroscience results. *Nature human behaviour*, 9(2):305–315, Nov 2025.
- [516] Zheheng Luo, Qianqian Xie, and Sophia Ananiadou. Citationsum: Citation-aware graph contrastive learning for scientific paper summarization. In *Proceedings of the ACM web conference 2023*, pages 1843–1852, Apr 2023.
- [517] Ziming Luo, Zonglin Yang, Zexin Xu, Wei Yang, and Xinya Du. Llm4sr: A survey on large language models for scientific research. *arXiv preprint arXiv:2501.04306*, 2025.
- [518] Andres M. Bran, Sam Cox, Oliver Schilter, Carlo Baldassari, Andrew D White, and Philippe Schwaller. Augmenting large language models with chemistry tools. *Nature Machine Intelligence*, 6(5):525–535, May 2024.
- [519] Pingchuan Ma, Tsun-Hsuan Wang, Minghao Guo, Zhiqing Sun, Joshua B Tenenbaum, Daniela Rus, Chuang Gan, and Wojciech Matusik. Llm and simulation as bilevel optimizers: A new paradigm to advance physical scientific discovery. *arXiv preprint arXiv:2405.09783*, 2024.
- [520] Qinyu Ma, Yuhao Zhou, and Jianfeng Li. Automated retrosynthesis planning of macromolecules using large language models and knowledge graphs. *Macromolecular Rapid Communications*, page 2500065, Feb 2025.
- [521] Shutian Ma, Heng Zhang, Chengzhi Zhang, and Xiaozhong Liu. Chronological citation recommendation with time preference. *Scientometrics*, 126:2991–3010, Feb 2021.
- [522] Yubo Ma, Zhibin Gou, Junheng Hao, Ruochen Xu, Shuohang Wang, Liangming Pan, Yujiu Yang, Yixin Cao, and Aixin Sun. SciAgent: Tool-augmented language models for scientific reasoning. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 15701–15736, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.880. URL <https://aclanthology.org/2024.emnlp-main.880/>.

- [523] Kojiro Machi, Seiji Akiyama, Yuuya Nagata, and Masaharu Yoshioka. A framework for reviewing the results of automated conversion of structured organic synthesis procedures from the literature. *Digital Discovery*, 4(1):172–180, Nov 2025.
- [524] Benjamin P MacLeod, Fraser GL Parlane, Thomas D Morrissey, Florian Häse, Loïc M Roch, Kevan E Dettelbach, Raphaell Moreira, Lars PE Yunker, Michael B Rooney, Joseph R Deeth, et al. Self-driving laboratory for accelerated discovery of thin-film materials. *Science Advances*, 6(20):eaaz8867, 2020.
- [525] Puja Maharjan. Benchmark for evaluation and analysis of citation recommendation models. *arXiv preprint arXiv:2412.07713*, 2024.
- [526] Bodhisattwa Prasad Majumder, Harshit Surana, Dhruv Agarwal, Sanchaita Hazra, Ashish Sabharwal, and Peter Clark. Position: data-driven discovery with large generative models. In *Forty-first International Conference on Machine Learning*, May 2024.
- [527] Bodhisattwa Prasad Majumder, Harshit Surana, Dhruv Agarwal, Bhavana Dalvi Mishra, Abhijeetsingh Meena, Aryan Prakhhar, Tirth Vora, Tushar Khot, Ashish Sabharwal, and Peter Clark. Discoverybench: Towards data-driven discovery with large language models. *arXiv preprint arXiv:2407.01725*, 2024.
- [528] Sepideh Maleki, Jan-Christian Huetter, Kangway V Chuang, David Richmond, Gabriele Scalia, and Tommaso Biancalani. Efficient fine-tuning of single-cell foundation models enables zero-shot molecular perturbation prediction. *arXiv preprint arXiv:2412.13478*, 2024.
- [529] Indrajeet Mandal, Jitendra Soni, Mohd Zaki, Morten M Smedskjaer, Katrin Wondraczek, Lothar Wondraczek, Nitya Nand Gosvami, and NM Krishnan. Autonomous microscopy experiments through large language model agents. *arXiv preprint arXiv:2501.10385*, 2024.
- [530] Benjamin S Manning, Kehang Zhu, and John J Horton. Automated social science: Language models as scientist and subjects. Technical report, National Bureau of Economic Research, Apr 2024.
- [531] Nacef Ben Mansour, Hamed Rahimi, and Motasem Alrahabi. How well do large language models extract keywords? a systematic evaluation on scientific corpora. In Peter Jansen, Bhavana Dalvi Mishra, Harsh Trivedi, Bodhisattwa Prasad Majumder, Tom Hope, Tushar Khot, Doug Downey, and Eric Horvitz, editors, *Proceedings of the 1st Workshop on AI and Scientific Discovery: Directions and Opportunities*, pages 13–21, Albuquerque, New Mexico, USA, May 2025. Association for Computational Linguistics. ISBN 979-8-89176-224-4. doi: 10.18653/v1/2025.aisd-main.2. URL <https://aclanthology.org/2025.aisd-main.2/>.
- [532] David M Markowitz. From complexity to clarity: How ai enhances perceptions of scientists and the public’s understanding of science. *PNAS nexus*, 3(9):pgae387, Sep 2024.
- [533] Anna Martin-Boyle, Aahan Tyagi, Marti A Hearst, and Dongyeop Kang. Shallow synthesis of knowledge in gpt-generated texts: A case study in automatic related work composition. *arXiv preprint arXiv:2402.12255*, 2024.
- [534] M Maruf, Arka Daw, Kazi Sajeed Mehrab, Harish Babu Manogaran, Abhilash Neog, Medha Sawhney, Mridul Khurana, James Balhoff, Yasin Bakis, Bahadir Altintas, et al. Vlm4bio: A benchmark dataset to evaluate pretrained vision-language models for trait discovery from biological images. *Advances in Neural Information Processing Systems*, 37:131035–131071, Sep 2024.

- [535] Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. ChartQA: A benchmark for question answering about charts with visual and logical reasoning. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio, editors, *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2263–2279, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.findings-acl.177. URL <https://aclanthology.org/2022.findings-acl.177/>.
- [536] Ahmed Masry, Mehrad Shahmohammadi, Md Rizwan Parvez, Enamul Hoque, and Shafiq Joty. ChartInstruct: Instruction tuning for chart comprehension and reasoning. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Findings of the Association for Computational Linguistics: ACL 2024*, pages 10387–10409, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.619. URL <https://aclanthology.org/2024.findings-acl.619/>.
- [537] Ahmed Masry, Megh Thakkar, Aayush Bajaj, Aaryaman Kartha, Enamul Hoque, and Shafiq Joty. ChartGemma: Visual instruction-tuning for chart reasoning in the wild. In Owen Rambow, Leo Wanner, Marianna Apidianaki, Hend Al-Khalifa, Barbara Di Eugenio, Steven Schockaert, Kareem Darwish, and Apoorv Agarwal, editors, *Proceedings of the 31st International Conference on Computational Linguistics: Industry Track*, pages 625–643, Abu Dhabi, UAE, January 2025. Association for Computational Linguistics. URL <https://aclanthology.org/2025.coling-industry.54/>.
- [538] Shray Mathur, Noah van der Vleuten, Kevin G Yager, and Esther Tsai. Vision: A modular ai assistant for natural human-instrument interaction at scientific user facilities. *Machine Learning: Science and Technology*, Jun 2025.
- [539] Blakeley B McShane, David Gal, and Adam Duhachek. Artificial intelligence and dichotomania. *Judgment and Decision Making*, 20:e23, Apr 2025.
- [540] Lennart Meincke, Karan Girotra, Gideon Nave, Christian Terwiesch, and Karl T. Ulrich. Using large language models for idea generation in innovation, Aug 2024. SSRN.
- [541] Lennart Meincke, Ethan R. Mollick, and Christian Terwiesch. Prompting diverse ideas: Increasing ai idea variance, Feb 2024. SSRN.
- [542] Chuizheng Meng, Sam Griesemer, Defu Cao, Sungyong Seo, and Yan Liu. When physics meets machine learning: A survey of physics-informed machine learning. *Machine Learning for Computational Science and Engineering*, 1(1):1–23, May 2025.
- [543] Fanqing Meng, Wenqi Shao, Quanfeng Lu, Peng Gao, Kaipeng Zhang, Yu Qiao, and Ping Luo. ChartAssistant: A universal chart multimodal language model via chart-to-table pre-training and multitask instruction tuning. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Findings of the Association for Computational Linguistics: ACL 2024*, pages 7775–7803, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.463. URL <https://aclanthology.org/2024.findings-acl.463/>.
- [544] Fanqing Meng, Wenqi Shao, Quanfeng Lu, Peng Gao, Kaipeng Zhang, Yu Qiao, and Ping Luo. Chartassistant: A universal chart multimodal language model via chart-to-table pre-training and multitask instruction tuning. *arXiv preprint arXiv:2401.02384*, 2024.
- [545] George Benneh Mensah. Artificial intelligence and ethics: a comprehensive review of bias mitigation, transparency, and accountability in ai systems. *Preprint, November*, 10(1), Nov 2023.

- [546] Amil Merchant, Simon Batzner, Samuel S Schoenholz, Muratahan Aykol, Gowoon Cheon, and Ekin Dogus Cubuk. Scaling deep learning for materials discovery. *Nature*, 624(7990):80–85, Nov 2023.
- [547] Dasha Metropolitansky and Jonathan Larson. Towards effective extraction and evaluation of factual claims. *arXiv preprint arXiv:2502.10855*, 2025.
- [548] Ning Miao, Yee Whye Teh, and Tom Rainforth. Selfcheck: Using llms to zero-shot check their own step-by-step reasoning. *arXiv preprint arXiv:2308.00436*, 2023.
- [549] Pawel Mieszczynek, Peter Corke, Courosh Mehanian, Paul D Dalton, and Dietmar W Huttmacher. Towards industry-ready additive manufacturing: Ai-enabled closed-loop control for 3d melt electrowriting. *Communications Engineering*, 3(1):158, 2024.
- [550] Xinyi Mou, Xuanwen Ding, Qi He, Liang Wang, Jingcong Liang, Xinnong Zhang, Libo Sun, Jiayu Lin, Jie Zhou, Xuanjing Huang, et al. From individual to society: A survey on social simulation driven by large language model-based agents. *arXiv preprint arXiv:2412.03563*, 2024.
- [551] Omar Moured, Sara Alzalabny, Anas Osman, Thorsten Schwarz, Karin Müller, and Rainer Stiefelwagen. Chartformer: A large vision language model for converting chart images into tactile accessible svgs. In *International Conference on Computers Helping People with Special Needs*, pages 299–305. Springer, May 2024.
- [552] Austin M Mroz, Annabel R Basford, Friedrich Hastedt, Isuru Shavindra Jayasekera, Irea Mosquera-Lois, Ruby Sedgwick, Pedro J Ballester, Joshua D Bocarsly, Ehecatl Antonio del Río Chanona, Matthew L Evans, et al. Cross-disciplinary perspectives on the potential for artificial intelligence across chemistry. *Chemical Society Reviews*, Apr 2025.
- [553] Panitan Muangkammuen, Fumiyo Fukumoto, Jiye Li, and Yoshimi Suzuki. Exploiting labeled and unlabeled data via transformer fine-tuning for peer-review score prediction. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 2233–2240, Dec 2022.
- [554] Arief Purnama Muharram and Ayu Purwarianti. Enhancing natural language inference performance with knowledge graph for covid-19 automated fact-checking in indonesian language. *arXiv preprint arXiv:2409.00061*, 2024.
- [555] Praveen Kumar Myakala, Anil Kumar Jonnalagadda, and Chiranjeevi Bura. Federated learning and data privacy: A review of challenges and opportunities. *International Journal of Research Publication and Reviews*, 5(12):10–55248, Jan 2024.
- [556] Sheshera Mysore, Debarati Das, Hancheng Cao, and Bahareh Sarrafzadeh. Prototypical human-ai collaboration behaviors from llm-assisted writing in the wild. *arXiv preprint arXiv:2505.16023*, 2025.
- [557] Inderjeet Nair, Jiaye Tan, Xiaotian Su, Anne Gere, Xu Wang, and Lu Wang. Closing the loop: Learning to generate writing feedback via language model simulated student revisions. *arXiv preprint arXiv:2410.08058*, 2024.
- [558] Siddharth Narayanan, James D Braza, Ryan-Rhys Griffiths, Manu Ponnampati, Albert Bou, Jon Laurent, Ori Kabeli, Geemi Wellawatte, Sam Cox, Samuel G Rodrigues, et al. Aviary: training language agents on challenging scientific tasks. *arXiv preprint arXiv:2412.21154*, 2024.

- [559] Deepak Nathani, Lovish Madaan, Nicholas Roberts, Nikolay Bashlykov, Ajay Menon, Vincent Moens, Amar Budhiraja, Despoina Magka, Vladislav Vorotilov, Gaurav Chaurasia, et al. Mlgym: A new framework and benchmark for advancing ai research agents. *arXiv preprint arXiv:2502.14499*, 2025.
- [560] Springer Nature. Snapp: Springer nature’s next-generation peer review system. <https://www.springernature.com/gp/snapp>, Dec 2023.
- [561] Vladimir Naumov, Diana Zagirova, Sha Lin, Yupeng Xie, Wenhao Gou, Anatoly Urban, Nina Tikhonova, Khadija Alawi, Mike Durymanov, Fedor Galkin, et al. Dora ai scientist: Multi-agent virtual research team for scientific exploration discovery and automated report generation. *bioRxiv*, Mar 2025.
- [562] Benjamin Newman, Yoonjoo Lee, Aakanksha Naik, Pao Siangliulue, Raymond Fok, Juho Kim, Daniel S Weld, Joseph Chee Chang, and Kyle Lo. ArxivDIGESTables: Synthesizing scientific literature into tables using language models. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 9612–9631, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.538. URL <https://aclanthology.org/2024.emnlp-main.538/>.
- [563] Izzy Newsham, Luka Kovačević, Richard Moulange, Nan Rosemary Ke, and Sach Mukherjee. Large language models for zero-shot inference of causal structures in biology. *arXiv preprint arXiv:2503.04347*, 2025.
- [564] Andy Nguyen, Yvonne Hong, Belle Dang, and Xiaoshan Huang. Human-ai collaboration patterns in ai-assisted academic writing. *Studies in Higher Education*, 49(5):847–864, Oct 2024.
- [565] Bo Ni and Markus J Buehler. Mechagents: Large language model multi-agent collaborations can solve mechanics problems, generate new data, and integrate knowledge. *Extreme Mechanics Letters*, 67: 102131, Mar 2024.
- [566] Ziqi Ni, Yahao Li, Kaijia Hu, Kunyuan Han, Ming Xu, Xingyu Chen, Fengqi Liu, Yicong Ye, and Shuxin Bai. Matpilot: an llm-enabled ai materials scientist under the framework of human-machine collaboration. *arXiv preprint arXiv:2411.08063*, 2024.
- [567] Harshit Nigam, Manasi Patwardhan, Lovekesh Vig, and Gautam Shroff. An interactive co-pilot for accelerated research ideation. In Su Lin Blodgett, Amanda Cercas Curry, Sunipa Dev, Michael Madaio, Ani Nenkova, Diyi Yang, and Ziang Xiao, editors, *Proceedings of the Third Workshop on Bridging Human–Computer Interaction and Natural Language Processing*, pages 60–73, Mexico City, Mexico, June 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.hcinlp-1.6. URL <https://aclanthology.org/2024.hcinlp-1.6/>.
- [568] Harshit Nigam, Manasi Patwardhan, Lovekesh Vig, and Gautam Shroff. Acceleron: A tool to accelerate research ideation. *arXiv preprint arXiv:2403.04382*, 2024.
- [569] Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan Wang, Yingbo Zhou, Silvio Savarese, and Caiming Xiong. Codegen: An open large language model for code with multi-turn program synthesis. *arXiv preprint arXiv:2203.13474*, 2022.
- [570] Kun-Peng Ning, Shuo Yang, Yuyang Liu, Jia-Yu Yao, Zhenhui Liu, Yonghong Tian, Yibing Song, and Li Yuan. PiCO: Peer review in LLMs based on consistency optimization. In *The Thirteenth International Conference on Learning Representations*, Jan 2025. URL <https://openreview.net/forum?id=sfQ6XpApfS>.

- [571] Kazuya Nishimura, Kuniaki Saito, Toshio Hirasawa, and Yoshitaka Ushiku. Toward related work generation with structure and novelty statement. In *Proceedings of the Fourth Workshop on Scholarly Document Processing (SDP 2024)*, pages 38–57, Aug 2024.
- [572] Haoyi Niu, Jianming Hu, Guyue Zhou, and Xianyuan Zhan. A comprehensive survey of cross-domain policy transfer for embodied agents. *arXiv preprint arXiv:2402.04580*, 2024.
- [573] Liang Niu, Nian Xue, and Christina Pöpper. Unveiling the sentinels: Assessing ai performance in cybersecurity peer review. *arXiv preprint arXiv:2309.05457*, 2023.
- [574] Alexander Novikov, Ngân Vu, Marvin Eisenberger, Emilien Dupont, Po-Sen Huang, Adam Zsolt Wagner, Sergey Shirobokov, Borislav Kozlovskii, Francisco JR Ruiz, Abbas Mehrabian, et al. Alphaevolve: A coding agent for scientific and algorithmic discovery. *Google DeepMind*, Jun 2025.
- [575] Charles O’Neill, Tirthankar Ghosal, Roberta Răileanu, Mike Walmsley, Thang Bui, Kevin Schawinski, and Ioana Ciucă. Sparks of science: Hypothesis generation using structured paper data. *arXiv preprint arXiv:2504.12976*, 2025.
- [576] OpenAI. Introducing deep research, Feb 2025. URL <https://openai.com/index/introducing-deep-research/>.
- [577] OpenAI. Gpt-4o search preview, May 2025. URL <https://platform.openai.com/docs/models/gpt-4o-search-preview>.
- [578] Raúl Ortega and José Manuel Gómez-Pérez. Sciclaims: An end-to-end generative system for biomedical claim analysis. *arXiv preprint arXiv:2503.18526*, 2025.
- [579] Zhinya Kawa Othman, Mohamed Mustaf Ahmed, Olalekan John Okesanya, Adamu Muhammad Ibrahim, Shuaibu Saidu Musa, Bryar A Hassan, Lanja Ibrahim Saeed, and Don Eliseo Lucero-Prisno III. Advancing drug discovery and development through gpt models: a review on challenges, innovations and future prospects. *Intelligence-Based Medicine*, page 100233, Mar 2025.
- [580] Colleen Ovelman, Shannon Kugley, Gerald Gartlehner, and Meera Viswanathan. The use of a large language model to create plain language summaries of evidence reviews in healthcare: A feasibility study. *Cochrane Evidence Synthesis and Methods*, 2(2):e12041, Feb 2024.
- [581] Ipek Ozkaya. Application of large language models to software engineering tasks: Opportunities, risks, and implications. *IEEE Software*, 40(3):4–8, Apr 2023.
- [582] Vishakh Padmakumar and He He. Machine-in-the-loop rewriting for creative image captioning. *arXiv preprint arXiv:2111.04193*, 2021.
- [583] Vishakh Padmakumar and He He. Does writing with language models reduce content diversity? *arXiv preprint arXiv:2309.05196*, 2023.
- [584] Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. Medmcqa: A large-scale multi-subject multi-choice dataset for medical domain question answering. In *Conference on health, inference, and learning*, pages 248–260. PMLR, Apr 2022.
- [585] Liangming Pan, Xiaobao Wu, Xinyuan Lu, Anh Tuan Luu, William Yang Wang, Min-Yen Kan, and Preslav Nakov. Fact-checking complex claims with program-guided reasoning. *arXiv preprint arXiv:2305.12744*, 2023.

- [586] Liangming Pan, Yunxiang Zhang, and Min-Yen Kan. Investigating zero-and few-shot generalization in fact verification. *arXiv preprint arXiv:2309.09444*, 2023.
- [587] Wenbo Pan, Zhichao Liu, Qiguang Chen, Xiangyang Zhou, Haining Yu, and Xiaohua Jia. The hidden dimensions of llm alignment: A multi-dimensional safety analysis. *arXiv preprint arXiv:2502.09674*, 2025.
- [588] Jing-Cheng Pang, Heng-Bo Fan, Pengyuan Wang, Jia-Hao Xiao, Nan Tang, Si-Hang Yang, Chengxing Jia, Sheng-Jun Huang, and Yang Yu. Empowering language models with active inquiry for deeper understanding. *arXiv preprint arXiv:2402.03719*, 2024.
- [589] Wei Pang, Kevin Qinghong Lin, Xiangru Jian, Xi He, and Philip Torr. Paper2poster: Towards multimodal poster automation from scientific papers. *arXiv preprint arXiv:2505.21497*, 2025.
- [590] Jong Inn Park, Maanas Taneja, Qianwen Wang, and Dongyeop Kang. Stealing creator’s workflow: A creator-inspired agentic framework with iterative feedback loop for improved scientific short-form generation. *arXiv preprint arXiv:2504.18805*, 2025.
- [591] Qizhi Pei, Lijun Wu, Kaiyuan Gao, Jinhua Zhu, Yue Wang, Zun Wang, Tao Qin, and Rui Yan. Leveraging biomolecule and natural language through multi-modal learning: A survey. *arXiv preprint arXiv:2403.01528*, 2024.
- [592] Vishnu S Pendyala, Karnavee Kamdar, and Kapil Mulchandani. Automated research review support using machine learning, large language models, and natural language processing. *Electronics*, 14(2): 256, Jan 2025.
- [593] Dengyun Peng, Yuhang Zhou, Qiguang Chen, Jinhao Liu, Jingjing Chen, and Libo Qin. Dlpo: Towards a robust, efficient, and generalizable prompt optimization framework from a deep-learning perspective. *arXiv preprint arXiv:2503.13413*, 2025.
- [594] Gal Peretz, Mousa Arraf, and Kira Radinsky. What if: Generating code to answer simulation questions in chemistry texts. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 1335–1344, Jul 2023.
- [595] Jérémy Perez, Corentin Léger, Marcela Ovando-Tellez, Chris Foulon, Joan Dussauld, Pierre-Yves Oudeyer, and Clément Moulin-Frier. Cultural evolution in populations of large language models. *arXiv preprint arXiv:2403.08882*, 2024.
- [596] Huy Nhat Phan, Tien N Nguyen, Phong X Nguyen, and Nghi DQ Bui. Hyperagent: Generalist software engineering agents to solve coding tasks at scale. *arXiv preprint arXiv:2409.16299*, 2024.
- [597] Iratxe Pinedo, Mikel Larrañaga, and Ana Arruarte. Arzigo: A recommendation system for scientific articles. *Information Systems*, 122:102367, May 2024.
- [598] Barbara Plank and Reinard van Dalen. Citetracked: A longitudinal dataset of peer reviews and citations. In *Proceedings of the 4th Joint Workshop on Bibliometric-enhanced Information Retrieval and Natural Language Processing for Digital Libraries (BIRNDL 2019) co-located with the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2019)*, pages 116–122. CEUR Workshop Proceedings (CEUR-WS. org), Jul 2019.
- [599] Stanislas Polu and Ilya Sutskever. Generative language modeling for automated theorem proving. *arXiv preprint arXiv:2009.03393*, 2020.

- [600] Dinesh K Pradhan, Joyita Chakraborty, Prasenjit Choudhary, and Subrata Nandi. An automated conflict of interest based greedy approach for conference paper assignment system. *Journal of Informetrics*, 14(2):101022, May 2020.
- [601] Carolina Pradier, Lucía Céspedes, and Vincent Larivière. A smack of all neighbouring languages: How multilingual is scholarly communication? *arXiv preprint arXiv:2504.21100*, 2025.
- [602] Shraman Pramanick, Rama Chellappa, and Subhashini Venugopalan. Spiqa: A dataset for multimodal question answering on scientific papers. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, Jul 2024. URL <https://openreview.net/forum?id=h3lddsY5nf>.
- [603] Adithya Pratapa and Teruko Mitamura. Estimating optimal context length for hybrid retrieval-augmented multi-document summarization. *arXiv preprint arXiv:2504.12972*, 2025.
- [604] S.V. Praveen, Pranshav Gajjar, Rajeev Kumar Ray, and Ashutosh Dutt. Crafting clarity: Leveraging large language models to decode consumer reviews. *Journal of Retailing and Consumer Services*, 81:103975, Nov 2024. ISSN 0969-6989. doi: <https://doi.org/10.1016/j.jretconser.2024.103975>. URL <https://www.sciencedirect.com/science/article/pii/S0969698924002716>.
- [605] Dongqi Pu and Vera Demberg. Rst-lora: A discourse-aware low-rank adaptation for long document abstractive summarization. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 2200–2220, May 2024.
- [606] Kevin Pu, KJ Kevin Feng, Tovi Grossman, Tom Hope, Bhavana Dalvi Mishra, Matt Latzke, Jonathan Bragg, Joseph Chee Chang, and Pao Siangliulue. Ideasynt: Iterative research idea development through evolving and composing idea facets with literature-grounded feedback. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*, pages 1–31, Apr 2025.
- [607] Yingming Pu, Tao Lin, and Hongyu Chen. Piflow: Principle-aware scientific discovery with multi-agent collaboration. *arXiv preprint arXiv:2505.15047*, 2025.
- [608] Sukannya Purkayastha, Zhuang Li, Anne Lauscher, Lizhen Qu, and Iryna Gurevych. Lazyreview a dataset for uncovering lazy thinking in nlp peer reviews. *arXiv preprint arXiv:2504.11042*, 2025.
- [609] I Made Putrama and Péter Martinek. Heterogeneous data integration: Challenges and opportunities. *Data in Brief*, page 110853, Oct 2024.
- [610] Edward O Pyzer-Knapp, Jed W Pitera, Peter WJ Staar, Seiji Takeda, Teodoro Laino, Daniel P Sanders, James Sexton, John R Smith, and Alessandro Curioni. Accelerating materials discovery using artificial intelligence, high performance computing and robotics. *npj Computational Materials*, 8(1):84, Apr 2022.
- [611] Chen Qian, Yufan Dang, Jiahao Li, Wei Liu, Zihao Xie, Yifei Wang, Weize Chen, Cheng Yang, Xin Cong, Xiaoyin Che, et al. Experiential co-learning of software-developing agents. *arXiv preprint arXiv:2312.17025*, 2023.
- [612] Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan Dang, Jiahao Li, Cheng Yang, Weize Chen, Yusheng Su, Xin Cong, et al. Chatdev: Communicative agents for software development. *arXiv preprint arXiv:2307.07924*, 2023.

- [613] Libo Qin, Xiao Xu, Wanxiang Che, Yue Zhang, and Ting Liu. Dynamic fusion network for multi-domain end-to-end task-oriented dialog. pages 6344–6354, July 2020. doi: 10.18653/v1/2020.acl-main.565. URL <https://aclanthology.org/2020.acl-main.565/>.
- [614] Libo Qin, Shijue Huang, Qiguang Chen, Chenran Cai, Yudi Zhang, Bin Liang, Wanxiang Che, and Ruifeng Xu. Mmsd2. 0: Towards a reliable multi-modal sarcasm detection system. *arXiv preprint arXiv:2307.07135*, 2023.
- [615] Libo Qin, Qiguang Chen, Hao Fei, Zhi Chen, Min Li, and Wanxiang Che. What factors affect multi-modal in-context learning? an in-depth exploration. *arXiv preprint arXiv:2410.20482*, 2024.
- [616] Libo Qin, Qiguang Chen, Xiachong Feng, Yang Wu, Yongheng Zhang, Yinghui Li, Min Li, Wanxiang Che, and Philip S Yu. Large language models meet nlp: A survey. *arXiv preprint arXiv:2405.12819*, 2024.
- [617] Libo Qin, Qiguang Chen, Yuhang Zhou, Zhi Chen, Yinghui Li, Lizi Liao, Min Li, Wanxiang Che, and Philip S Yu. A survey of multilingual large language models. *Patterns*, 6(1), Jan 2025. URL [https://www.cell.com/patterns/fulltext/S2666-3899\(24\)00290-3](https://www.cell.com/patterns/fulltext/S2666-3899(24)00290-3).
- [618] Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, et al. Toolllm: Facilitating large language models to master 16000+ real-world apis. *arXiv preprint arXiv:2307.16789*, 2023.
- [619] 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.
- [620] Yansheng Qiu, Haoquan Zhang, Zhaopan Xu, Ming Li, Diping Song, Zheng Wang, and Kaipeng Zhang. Ai idea bench 2025: Ai research idea generation benchmark. *arXiv preprint arXiv:2504.14191*, 2025.
- [621] Chen Qu, Liu Yang, Cen Chen, Minghui Qiu, W Bruce Croft, and Mohit Iyyer. Open-retrieval conversational question answering. In *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval*, pages 539–548, Jul 2020.
- [622] Yuxiao Qu, Tianjun Zhang, Naman Garg, and Aviral Kumar. Recursive introspection: Teaching language model agents how to self-improve. *Advances in Neural Information Processing Systems*, 37: 55249–55285, Dec 2024.
- [623] Marissa Radensky, Simra Shahid, Raymond Fok, Pao Siangliulue, Tom Hope, and Daniel S Weld. Scideator: Human-llm scientific idea generation grounded in research-paper facet recombination. *arXiv preprint arXiv:2409.14634*, 2024.
- [624] Ilija Radosavovic, Tete Xiao, Bike Zhang, Trevor Darrell, Jitendra Malik, and Koushil Sreenath. Real-world humanoid locomotion with reinforcement learning. *Science Robotics*, 9(89):eadi9579, 2024.
- [625] Mohammad Raeini. On the rise of new mathematical spaces and towards ai-driven scientific discovery. *Available at SSRN*, Mar 2025.
- [626] Aditi Raghunathan, Nihar B Shah, et al. Vulnerability of text-matching in ml/ai conference reviewer assignments to collusions. *arXiv preprint arXiv:2412.06606*, 2024.

- [627] Daking Rai, Yilun Zhou, Shi Feng, Abulhair Saparov, and Ziyu Yao. A practical review of mechanistic interpretability for transformer-based language models. *arXiv preprint arXiv:2407.02646*, 2024.
- [628] Maziar Raissi, Paris Perdikaris, and George E Karniadakis. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational physics*, 378:686–707, Feb 2019.
- [629] Maziar Raissi, Paris Perdikaris, Nazanin Ahmadi, and George Em Karniadakis. Physics-informed neural networks and extensions. *arXiv preprint arXiv:2408.16806*, 2024.
- [630] Mahyar Rajabi-Kochi, Negareh Mahboubi, Aseem Partap Singh Gill, and Seyed Mohamad Moosavi. Adaptive representation of molecules and materials in bayesian optimization. *Chemical Science*, 16 (13):5464–5474, Feb 2025.
- [631] Sriram Ranga, Rui Mao, Erik Cambria, and Anupam Chattopadhyay. The plagiarism singularity conjecture. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 10245–10255, Apr 2025.
- [632] Delip Rao, Jonathan Young, Thomas Dietterich, and Chris Callison-Burch. Withdrarxiv: A large-scale dataset for retraction study. *arXiv preprint arXiv:2412.03775*, 2024.
- [633] Delip Rao, Weiqiu You, Eric Wong, and Chris Callison-Burch. Nsf-scify: Mining the nsf awards database for scientific claims. *arXiv preprint arXiv:2503.08600*, 2025.
- [634] Vishisht Rao, Aounon Kumar, Himabindu Lakkaraju, and Nihar B Shah. Detecting llm-written peer reviews. *arXiv preprint arXiv:2503.15772*, 2025.
- [635] Md Mahinur Rashid, Hasin Kawsar Jahan, Annysha Huzzat, Riyasaat Ahmed Rahul, Tamim Bin Zakir, Farhana Meem, Md Saddam Hossain Mukta, and Swakkhar Shatabda. Text2chart: A multi-staged chart generator from natural language text. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, pages 3–16. Springer, May 2022.
- [636] Shaina Raza, Brian Schwartz, and Laura C Rosella. Coquad: a covid-19 question answering dataset system, facilitating research, benchmarking, and practice. *BMC bioinformatics*, 23(1):210, Jun 2022.
- [637] Chandan K Reddy and Parshin Shojaee. Towards scientific discovery with generative ai: Progress, opportunities, and challenges. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 28601–28609, Jan 2025.
- [638] Tohida Rehman, Debarshi Kumar Sanyal, and Samiran Chattopadhyay. Can pre-trained language models generate titles for research papers? In *International Conference on Asian Digital Libraries*, pages 154–170. Springer, Dec 2025.
- [639] Shuo Ren, Pu Jian, Zhenjiang Ren, Chunlin Leng, Can Xie, and Jiajun Zhang. Towards scientific intelligence: A survey of llm-based scientific agents. *arXiv preprint arXiv:2503.24047*, 2025.
- [640] Jonathan Roberts, Kai Han, Neil Houlsby, and Samuel Albanie. Scifibench: Benchmarking large multimodal models for scientific figure interpretation. *arXiv preprint arXiv:2405.08807*, 2024.
- [641] Zachary Robertson. Gpt4 is slightly helpful for peer-review assistance: A pilot study. *arXiv preprint arXiv:2307.05492*, 2023.

- [642] Ambrose Robinson, William Thorne, Ben P Wu, Abdullah Pandor, Munira Essat, Mark Stevenson, and Xingyi Song. Bio-sieve: exploring instruction tuning large language models for systematic review automation. *arXiv preprint arXiv:2308.06610*, 2023.
- [643] Juan A Rodriguez, David Vazquez, Issam Laradji, Marco Pedersoli, and Pau Rodriguez. Figgen: Text to scientific figure generation. *arXiv preprint arXiv:2306.00800*, 2023.
- [644] Juan A Rodriguez, Xiangru Jian, Siba Smarak Panigrahi, Tianyu Zhang, Aarash Feizi, Abhay Puri, Akshay Kalkunte Suresh, François Savard, Ahmed Masry, Shravan Nayak, et al. Bigdocs: An open dataset for training multimodal models on document and code tasks. In *The Thirteenth International Conference on Learning Representations*, Jan 2025.
- [645] Juan A Rodriguez, Abhay Puri, Shubham Agarwal, Issam H Laradji, Pau Rodriguez, Sai Rajeswar, David Vazquez, Christopher Pal, and Marco Pedersoli. Starvector: Generating scalable vector graphics code from images and text. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 16175–16186, Apr 2025.
- [646] D Gordon Rohman. Pre-writing: The stage of discovery in the writing process. *College Composition & Communication*, 16(2):106–112, May 1965.
- [647] Jennifer E Rood, Anna Hupalowska, and Aviv Regev. Toward a foundation model of causal cell and tissue biology with a perturbation cell and tissue atlas. *Cell*, 187(17):4520–4545, Aug 2024.
- [648] Yusuf Roohani, Andrew Lee, Qian Huang, Jian Vora, Zachary Steinhart, Kexin Huang, Alexander Marson, Percy Liang, and Jure Leskovec. Biodiscoveryagent: An ai agent for designing genetic perturbation experiments. *arXiv preprint arXiv:2405.17631*, 2024.
- [649] Emely Rosbach, Jonathan Ganz, Jonas Ammeling, Andreas Riener, and Marc Aubreville. Automation bias in ai-assisted medical decision-making under time pressure in computational pathology. In *BVM Workshop*, pages 129–134. Springer, Mar 2025.
- [650] Zhyar Rzgar K Rostam and Gábor Kertész. Fine-tuning large language models for scientific text classification: A comparative study. In *2024 IEEE 6th International Symposium on Logistics and Industrial Informatics (LINDI)*, pages 000233–000238. IEEE, Oct 2024.
- [651] Tony Rousmaniere, Xu Li, Yimeng Zhang, and Siddharth Shah. Large language models as mental health resources: Patterns of use in the united states, Mar 2025.
- [652] Pritam Roy and Dhananjay Datta. Ai-driven discovery: The transformative impact of machine learning on research and development. In *Evolving Landscapes of Research and Development: Trends, Challenges, and Opportunities*, pages 29–52. IGI Global Scientific Publishing, Jun 2025.
- [653] Sayar Ghosh Roy and Jiawei Han. Ilciter: Evidence-grounded interpretable local citation recommendation. *arXiv preprint arXiv:2403.08737*, 2024.
- [654] Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, et al. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*, 2023.
- [655] Kai Ruan, Xuan Wang, Jixiang Hong, and Hao Sun. Liveideabench: Evaluating llms’ scientific creativity and idea generation with minimal context. *arXiv preprint arXiv:2412.17596*, 2024.

- [656] Yixiang Ruan, Chenyin Lu, Ning Xu, Yuchen He, Yixin Chen, Jian Zhang, Jun Xuan, Jianzhang Pan, Qun Fang, Hanyu Gao, et al. An automatic end-to-end chemical synthesis development platform powered by large language models. *Nature communications*, 15(1):10160, Nov 2024.
- [657] Masoud Jalili Sabet, Philipp Dufter, François Yvon, and Hinrich Schütze. Simalign: High quality word alignments without parallel training data using static and contextualized embeddings. *arXiv preprint arXiv:2004.08728*, 2020.
- [658] Henrik Skaug Sætra. The rise of the research automaton: Science as process or product in the era of generative ai? *Available at SSRN 5219722*, Apr 2025.
- [659] Tanik Saikh, Tirthankar Ghosal, Amish Mittal, Asif Ekbal, and Pushpak Bhattacharyya. Scienceqa: A novel resource for question answering on scholarly articles. *International Journal on Digital Libraries*, 23(3):289–301, Jul 2022.
- [660] Shubhra Kanti Karmaker Santu, Sanjeev Kumar Sinha, Naman Bansal, Alex Knipper, Souvika Sarkar, John Salvador, Yash Mahajan, Sri Guttikonda, Mousumi Akter, Matthew Freestone, et al. Prompting llms to compose meta-review drafts from peer-review narratives of scholarly manuscripts. *arXiv preprint arXiv:2402.15589*, 2024.
- [661] Aishik Sanyal, Samuel Schapiro, Sumuk Shashidhar, Royce Moon, Lav R Varshney, and Dilek Hakkani-Tur. Spark: A system for scientifically creative idea generation. *arXiv preprint arXiv:2504.20090*, 2025.
- [662] Burcu Sayin, Ipek Baris Schlicht, Ngoc Vo Hong, Sara Allievi, Jacopo Staiano, Pasquale Minervini, and Andrea Passerini. Medsyn: Enhancing diagnostics with human-ai collaboration. *arXiv preprint arXiv:2506.14774*, 2025.
- [663] Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36:68539–68551, Dec 2023.
- [664] Emma Schleiger, Claire Mason, Claire Naughtin, Andrew Reeson, and Cecile Paris. Collaborative intelligence: A scoping review of current applications. *Applied Artificial Intelligence*, 38(1):2327890, Mar 2024.
- [665] Samuel Schmidgall and Michael Moor. Agentrxiv: Towards collaborative autonomous research. *arXiv preprint arXiv:2503.18102*, 2025.
- [666] Samuel Schmidgall, Yusheng Su, Ze Wang, Ximeng Sun, Jialian Wu, Xiaodong Yu, Jiang Liu, Zicheng Liu, and Emad Barsoum. Agent laboratory: Using llm agents as research assistants. *arXiv preprint arXiv:2501.04227*, 2025.
- [667] Antony Seabra, Claudio Cavalcante, Joao Nepomuceno, Lucas Lago, Nicolaas Ruberg, and Sergio Lifschitz. Dynamic multi-agent orchestration and retrieval for multi-source question-answer systems using large language models. *arXiv preprint arXiv:2412.17964*, 2024.
- [668] Andrew W Senior, Richard Evans, John Jumper, James Kirkpatrick, Laurent Sifre, Tim Green, Chongli Qin, Augustin Židek, Alexander WR Nelson, Alex Bridgland, et al. Improved protein structure prediction using potentials from deep learning. *Nature*, 577(7792):706–710, Jan 2020.

- [669] Junpyo Seo, Dongwan Kim, Jaewook Jeong, Inkyu Park, and Junho Min. FlavorDiffusion: Modeling food-chemical interactions with diffusion. In Peter Jansen, Bhavana Dalvi Mishra, Harsh Trivedi, Bodhisattwa Prasad Majumder, Tom Hope, Tushar Khot, Doug Downey, and Eric Horvitz, editors, *Proceedings of the 1st Workshop on AI and Scientific Discovery: Directions and Opportunities*, pages 70–77, Albuquerque, New Mexico, USA, May 2025. Association for Computational Linguistics. ISBN 979-8-89176-224-4. doi: 10.18653/v1/2025.aisd-main.7. URL <https://aclanthology.org/2025.aisd-main.7/>.
- [670] Sergio A Serrano, Jose Martinez-Carranza, and L Enrique Sucar. Knowledge transfer for cross-domain reinforcement learning: a systematic review. *IEEE Access*, Jul 2024.
- [671] Abdul Shahid, Muhammad Tanvir Afzal, Moloud Abdar, Mohammad Ehsan Basiri, Xujuan Zhou, Neil Y Yen, and Jia-Wei Chang. Insights into relevant knowledge extraction techniques: a comprehensive review. *The Journal of Supercomputing*, 76:1695–1733, Oct 2020.
- [672] Mohamed H Shahin, Srijib Goswami, Sebastian Lobentanzer, and Brian W Corrigan. Agents for change: Artificial intelligent workflows for quantitative clinical pharmacology and translational sciences. *Clinical and Translational Science*, 18(3):e70188, Mar 2025.
- [673] Yijia Shao, Yucheng Jiang, Theodore A Kanell, Peter Xu, Omar Khattab, and Monica S Lam. Asisting in writing wikipedia-like articles from scratch with large language models. *arXiv preprint arXiv:2402.14207*, 2024.
- [674] Zachary Sheldon and Peeyush Kumar. Economic anthropology in the era of generative artificial intelligence. *arXiv preprint arXiv:2410.15238*, 2024.
- [675] Xiang Shi, Jiawei Liu, Yinpeng Liu, Qikai Cheng, and Wei Lu. Every part matters: Integrity verification of scientific figures based on multimodal large language models. *arXiv preprint arXiv:2407.18626*, 2024.
- [676] Zhengliang Shi, Shen Gao, Zhen Zhang, Xiuying Chen, Zhumin Chen, Pengjie Ren, and Zhaochun Ren. Towards a unified framework for reference retrieval and related work generation. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 5785–5799, Dec 2023.
- [677] Zhengliang Shi, Lingyong Yan, Dawei Yin, Suzan Verberne, Maarten de Rijke, and Zhaochun Ren. Iterative self-incentivization empowers large language models as agentic searchers. *arXiv preprint arXiv:2505.20128*, 2025.
- [678] Hyungyu Shin, Jingyu Tang, Yoonjoo Lee, Nayoung Kim, Hyunseung Lim, Ji Yong Cho, Hwajung Hong, Moontae Lee, and Juho Kim. Automatically evaluating the paper reviewing capability of large language models. *arXiv preprint arXiv:2502.17086*, 2025.
- [679] Hyungyu Shin, Jingyu Tang, Yoonjoo Lee, Nayoung Kim, Hyunseung Lim, Ji Yong Cho, Hwajung Hong, Moontae Lee, and Juho Kim. Mind the blind spots: A focus-level evaluation framework for llm reviews. *arXiv preprint arXiv:2502.17086*, 2025.
- [680] Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems*, 36:8634–8652, Dec 2023.

- [681] Parshin Shojaee, Kazem Meidani, Shashank Gupta, Amir Barati Farimani, and Chandan K. Reddy. LLM-SR: Scientific equation discovery via programming with large language models. In *The Thirteenth International Conference on Learning Representations*, Jan 2025. URL <https://openreview.net/forum?id=m2nmp8P5in>.
- [682] Parshin Shojaee, Ngoc-Hieu Nguyen, Kazem Meidani, Amir Barati Farimani, Khoa D Doan, and Chandan K Reddy. Llm-srbench: A new benchmark for scientific equation discovery with large language models. *arXiv preprint arXiv:2504.10415*, 2025.
- [683] Parshin Shojaee, Ngoc-Hieu Nguyen, Kazem Meidani, Amir Barati Farimani, Khoa D Doan, and Chandan K. Reddy. LLM-SRBench: A new benchmark for scientific equation discovery with large language models. In *Forty-second International Conference on Machine Learning*, May 2025. URL <https://openreview.net/forum?id=SyQPizJVWY>.
- [684] Chenglei Si, Diyi Yang, and Tatsunori Hashimoto. Can llms generate novel research ideas? a large-scale human study with 100+ nlp researchers. *arXiv preprint arXiv:2409.04109*, 2024.
- [685] Shoaib Ahmed Siddiqui, Yanzhi Chen, Juyeon Heo, Menglin Xia, and Adrian Weller. On evaluating LLMs’ capabilities as functional approximators: A Bayesian evaluation framework. In Owen Rambow, Leo Wanner, Marianna Apidianaki, Hend Al-Khalifa, Barbara Di Eugenio, and Steven Schockaert, editors, *Proceedings of the 31st International Conference on Computational Linguistics*, pages 5826–5835, Abu Dhabi, UAE, January 2025. Association for Computational Linguistics. URL <https://aclanthology.org/2025.coling-main.388/>.
- [686] Aditi Singh, Abul Ehtesham, Saket Kumar, and Tala Talaei Khoei. Agentic retrieval-augmented generation: A survey on agentic rag. *arXiv preprint arXiv:2501.09136*, 2025.
- [687] Amanpreet Singh, Joseph Chee Chang, Chloe Anastasiades, Dany Haddad, Aakanksha Naik, Amber Tanaka, Angele Zamarron, Cecile Nguyen, Jena D Hwang, Jason Dunkleberger, et al. Ai2 scholar qa: Organized literature synthesis with attribution. *arXiv preprint arXiv:2504.10861*, 2025.
- [688] Ashish Singh, Prateek Agarwal, Zixuan Huang, Arpita Singh, Tong Yu, Sungchul Kim, Victor Bursztyn, Nikos Vlassis, and Ryan A Rossi. Figcaps-hf: A figure-to-caption generative framework and benchmark with human feedback. *arXiv preprint arXiv:2307.10867*, 2023.
- [689] Sanjay Singh and Amaresh Chakrabarti. Supporting assessment of novelty of design problems using concept of problem sapphire. *arXiv preprint arXiv:2410.18629*, 2024.
- [690] Shruti Singh, Mayank Singh, and Pawan Goyal. Compare: a taxonomy and dataset of comparison discussions in peer reviews. In *2021 ACM/IEEE Joint Conference on Digital Libraries (JCDL)*, pages 238–241. IEEE, Aug 2021.
- [691] Shruti Singh, Nandan Sarkar, and Arman Cohan. Scidqa: A deep reading comprehension dataset over scientific papers. *arXiv preprint arXiv:2411.05338*, 2024.
- [692] Karan Singhal, Shekoofeh Azizi, Tao Tu, S. Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, Perry Payne, Martin Seneviratne, Paul Gamble, Chris Kelly, Abubakr Babiker, Nathanael Schärli, Aakanksha Chowdhery, Philip Mansfield, Dina Demner-Fushman, Blaise Agüera y Arcas, Dale Webster, Greg S. Corrado, Yossi Matias, Katherine Chou, Juraj Gottweis, Nenad Tomasev, Yun Liu, Alvin Rajkomar, Joelle Barral, Christopher Semturs, Alan Karthikesalingam, and Vivek Natarajan. Large language models encode clinical knowledge.

Nature, 620(7972):172–180, Aug 2023. ISSN 1476-4687. doi: 10.1038/s41586-023-06291-2. URL <https://doi.org/10.1038/s41586-023-06291-2>.

- [693] Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Mohamed Amin, Le Hou, Kevin Clark, Stephen R Pfohl, Heather Cole-Lewis, et al. Toward expert-level medical question answering with large language models. *Nature Medicine*, pages 1–8, Jan 2025.
- [694] Arnab Sinha, Zhihong Shen, Yang Song, Hao Ma, Darrin Eide, Bo-June Hsu, and Kuansan Wang. An overview of microsoft academic service (mas) and applications. In *Proceedings of the 24th international conference on world wide web*, pages 243–246, May 2015.
- [695] Shiven Sinha, Shashwat Goel, Ponnurangam Kumaraguru, Jonas Geiping, Matthias Bethge, and Ameya Prabhu. Can language models falsify? evaluating algorithmic reasoning with counterexample creation. *arXiv preprint arXiv:2502.19414*, 2025.
- [696] Michael D Skarlinski, Sam Cox, Jon M Laurent, James D Braza, Michaela Hinks, Michael J Hammerling, Manvitha Ponnampati, Samuel G Rodriques, and Andrew D White. Language agents achieve superhuman synthesis of scientific knowledge. *arXiv preprint arXiv:2409.13740*, 2024.
- [697] Khachik Smbatyan, Tsolak Ghukasyan, Tigran Aghajanyan, Hovhannes Dabaghyan, Sergey Adamyan, Aram Bughdaryan, Vahagn Altunyan, Gagik Navasardyan, Aram Davtyan, Anush Hakobyan, et al. Can ai agents design and implement drug discovery pipelines? *arXiv preprint arXiv:2504.19912*, 2025.
- [698] IEEE Computer Society. How to make peer review recommendations and decisions. <https://www.computer.org/publications/making-peer-review-recommendations>.
- [699] Gleb Vitalevich Solovev, Alina Borisovna Zhidkovskaya, Anastasia Orlova, Anastasia Vepreva, Tonkii Ilya, Rodion Golovinskii, Nina Gubina, Denis Chistiakov, Timur A Aliev, Ivan Poddiakov, et al. Towards llm-driven multi-agent pipeline for drug discovery: Neurodegenerative diseases case study. In *2nd AI4Research Workshop: Towards a Knowledge-grounded Scientific Research Lifecycle*, Dec 2024.
- [700] Guijin Son, Jiwoo Hong, Honglu Fan, Heejeong Nam, Hyunwoo Ko, Seungwon Lim, Jinyeop Song, Jinha Choi, Gonçalo Paulo, Youngjae Yu, et al. When ai co-scientists fail: Spot-a benchmark for automated verification of scientific research. *arXiv preprint arXiv:2505.11855*, 2025.
- [701] Cuiping Song and Yanping Song. Enhancing academic writing skills and motivation: assessing the efficacy of chatgpt in ai-assisted language learning for efl students. *Frontiers in Psychology*, 14: 1260843, Dec 2023.
- [702] Jinwang Song, Yanxin Song, Guangyu Zhou, Wenhui Fu, Kunli Zhang, and Hongying Zan. Enhancing chinese essay discourse logic evaluation through optimized fine-tuning of large language models. In *CCF International Conference on Natural Language Processing and Chinese Computing*, pages 342–352. Springer, Nov 2024.
- [703] Peiyang Song, Kaiyu Yang, and Anima Anandkumar. Towards large language models as copilots for theorem proving in lean. *arXiv preprint arXiv:2404.12534*, 2024.
- [704] Tao Song, Man Luo, Xiaolong Zhang, Linjiang Chen, Yan Huang, Jiaqi Cao, Qing Zhu, Daobin Liu, Baicheng Zhang, Gang Zou, et al. A multiagent-driven robotic ai chemist enabling autonomous chemical research on demand. *Journal of the American Chemical Society*, 147(15):12534–12545, Mar 2025.

- [705] Yuda Song, Hanlin Zhang, Carson Eisenach, Sham Kakade, Dean Foster, and Udaya Ghai. Mind the gap: Examining the self-improvement capabilities of large language models. *arXiv preprint arXiv:2412.02674*, 2024.
- [706] Zhilong Song, Minggang Ju, Chunjin Ren, Qiang Li, Chongyi Li, Qionghua Zhou, and Jinlan Wang. Llm-feynman: Leveraging large language models for universal scientific formula and theory discovery. *arXiv preprint arXiv:2503.06512*, 2025.
- [707] Evan Walter Clark Spotte-Smith. Considering the ethics of large machine learning models in the chemical sciences. Mar 2025.
- [708] Henry W Sprueill, Carl Edwards, Khushbu Agarwal, Mariefel V Olarte, Udishnu Sanyal, Conrad Johnston, Hongbin Liu, Heng Ji, and Sutanay Choudhury. Chemreasoner: Heuristic search over a large language model’s knowledge space using quantum-chemical feedback. *arXiv preprint arXiv:2402.10980*, 2024.
- [709] Liana Spytka. The use of artificial intelligence in psychotherapy: development of intelligent therapeutic systems. *BMC psychology*, 13(1):175, Feb 2025.
- [710] Karthik Sreedhar, Alice Cai, Jenny Ma, Jeffrey V Nickerson, and Lydia B Chilton. Simulating cooperative prosocial behavior with multi-agent llms: Evidence and mechanisms for ai agents to inform policy decisions. In *Proceedings of the 30th International Conference on Intelligent User Interfaces*, pages 1272–1286, Mar 2025.
- [711] Kaustubh Sridhar, Souradeep Dutta, Dinesh Jayaraman, and Insup Lee. Regent: A retrieval-augmented generalist agent that can act in-context in new environments. In *NeurIPS 2024 Workshop on Open-World Agents*, Dec 2024.
- [712] Moritz Staudinger, Wojciech Kusa, Florina Piroi, and Allan Hanbury. An analysis of tasks and datasets in peer reviewing. In *Proceedings of the Fourth Workshop on Scholarly Document Processing (SDP 2024)*, pages 257–268, Aug 2024.
- [713] Vaio Stergiopoulos, Michael Vassilakopoulos, Eleni Tousidou, and Antonio Corral. An academic recommender system on large citation data based on clustering, graph modeling and deep learning. *Knowledge and Information Systems*, 66(8):4463–4496, Apr 2024.
- [714] Noy Sternlicht and Tom Hope. Chimera: A knowledge base of idea recombination in scientific literature. *arXiv preprint arXiv:2505.20779*, 2025.
- [715] Jonathan M Stokes, Kevin Yang, Kyle Swanson, Wengong Jin, Andres Cubillos-Ruiz, Nina M Donghia, Craig R MacNair, Shawn French, Lindsey A Carfrae, Zohar Bloom-Ackermann, et al. A deep learning approach to antibiotic discovery. *Cell*, 180(4):688–702, Feb 2020.
- [716] James W. A. Strachan, Dalila Albergo, Giulia Borghini, Oriana Pansardi, Eugenio Scaliti, Saurabh Gupta, Krati Saxena, Alessandro Rufo, Stefano Panzeri, Guido Manzi, Michael S. A. Graziano, and Cristina Becchio. Testing theory of mind in large language models and humans. *Nature Human Behaviour*, 8:1285–1295, May 2024. doi: 10.1038/s41562-024-01895-1. URL <https://www.nature.com/articles/s41562-024-01895-1>.
- [717] 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.

- [718] Yuan Sui, Mengyu Zhou, Mingjie Zhou, Shi Han, and Dongmei Zhang. Table meets llm: Can large language models understand structured table data? a benchmark and empirical study. In *Proceedings of the 17th ACM International Conference on Web Search and Data Mining, WSDM '24*, page 645–654, New York, NY, USA, Mar 2024. Association for Computing Machinery. ISBN 9798400703713. doi: 10.1145/3616855.3635752. URL <https://doi.org/10.1145/3616855.3635752>.
- [719] Yuan Sui, Mengyu Zhou, Mingjie Zhou, Shi Han, and Dongmei Zhang. Table meets llm: Can large language models understand structured table data? a benchmark and empirical study. In *Proceedings of the 17th ACM International Conference on Web Search and Data Mining*, pages 645–654, Mar 2024.
- [720] Purin Sukpanichnant, Anna Rapberger, and Francesca Toni. Peerarg: Argumentative peer review with llms. *arXiv preprint arXiv:2409.16813*, 2024.
- [721] Bo Sun, Baoxin Wang, Yixuan Wang, Wanxiang Che, Dayong Wu, Shijin Wang, and Ting Liu. Csed: A chinese semantic error diagnosis corpus. *arXiv preprint arXiv:2305.05183*, 2023.
- [722] Chongyan Sun, Ken Lin, Shiwei Wang, Hulong Wu, Chengfei Fu, and Zhen Wang. Lalaeval: A holistic human evaluation framework for domain-specific large language models. In *First Conference on Language Modeling*, Aug 2024.
- [723] Hao Sun, Yunyi Shen, and Mihaela van der Schaar. Openreview should be protected and leveraged as a community asset for research in the era of large language models. *arXiv preprint arXiv:2505.21537*, 2025.
- [724] Jiankai Sun, Chuanyang Zheng, Enze Xie, Zhengying Liu, Ruihang Chu, Jianing Qiu, Jiaqi Xu, Mingyu Ding, Hongyang Li, Mengzhe Geng, et al. A survey of reasoning with foundation models. *arXiv preprint arXiv:2312.11562*, 2023.
- [725] Lu Sun, Stone Tao, Junjie Hu, and Steven P Dow. Metawriter: Exploring the potential and perils of ai writing support in scientific peer review. *Proceedings of the ACM on Human-Computer Interaction*, 8 (CSCW1):1–32, Apr 2024.
- [726] Qiang Sun, Yuanyi Luo, Wenxiao Zhang, Sirui Li, Jichunyang Li, Kai Niu, Xiangrui Kong, and Wei Liu. Docs2kg: A human-llm collaborative approach to unified knowledge graph construction from heterogeneous documents. In *Companion Proceedings of the ACM on Web Conference 2025*, pages 801–804, May 2025.
- [727] Qiushi Sun, Zhoumianze Liu, Chang Ma, Zichen Ding, Fangzhi Xu, Zhangyue Yin, Haiteng Zhao, Zhenyu Wu, Kanzhi Cheng, Zhaoyang Liu, et al. Scienceboard: Evaluating multimodal autonomous agents in realistic scientific workflows. *arXiv preprint arXiv:2505.19897*, 2025.
- [728] Tao Sun, Enhao Pan, Zhengkai Yang, Kaixin Sui, Jiajun Shi, Xianfu Cheng, Tongliang Li, Wenhao Huang, Ge Zhang, Jian Yang, et al. P2p: Automated paper-to-poster generation and fine-grained benchmark. *arXiv preprint arXiv:2505.17104*, 2025.
- [729] Jayaprakash Sundararaj, Akhil Vyas, and Benjamin Gonzalez-Maldonado. Automated latex code generation from handwritten math expressions using vision transformer. *arXiv preprint arXiv:2412.03853*, 2024.
- [730] Gaurav Suri, Lily R Slater, Ali Ziaee, and Morgan Nguyen. Do large language models show decision heuristics similar to humans? a case study using gpt-3.5. *Journal of Experimental Psychology: General*, 153(4):1066, May 2024.

- [731] Teo Susnjak, Peter Hwang, Napoleon Reyes, Andre LC Barczak, Timothy McIntosh, and Surangika Ranathunga. Automating research synthesis with domain-specific large language model fine-tuning. *ACM Transactions on Knowledge Discovery from Data*, 19(3):1–39, Mar 2025.
- [732] Kyle Swanson, Wesley Wu, Nash L Bulaong, John E Pak, and James Zou. The virtual lab: Ai agents design new sars-cov-2 nanobodies with experimental validation. *bioRxiv*, pages 2024–11, Nov 2024.
- [733] Jaroslaw Szumega, Lamine Bougueroua, Blerina Gkotse, Pierre Jouvelot, and Federico Ravotti. The open review-based (orb) dataset: Towards automatic assessment of scientific papers and experiment proposals in high-energy physics. *arXiv preprint arXiv:2312.04576*, 2023.
- [734] Nathan J Szymanski, Bernardus Rendy, Yuxing Fei, Rishi E Kumar, Tanjin He, David Milsted, Matthew J McDermott, Max Gallant, Ekin Dogus Cubuk, Amil Merchant, et al. An autonomous laboratory for the accelerated synthesis of novel materials. *Nature*, 624(7990):86–91, Nov 2023.
- [735] Pawin Taechoyotin and Daniel Acuna. Remor: Automated peer review generation with llm reasoning and multi-objective reinforcement learning. *arXiv preprint arXiv:2505.11718*, 2025.
- [736] Pawin Taechoyotin, Guanchao Wang, Tong Zeng, Bradley Sides, and Daniel Acuna. Mamorx: Multi-agent multi-modal scientific review generation with external knowledge. In *Neurips 2024 Workshop Foundation Models for Science: Progress, Opportunities, and Challenges*, Oct 2024.
- [737] Iman Tahamtan, Askar Safipour Afshar, and Khadijeh Ahamdzadeh. Factors affecting number of citations: a comprehensive review of the literature. *Scientometrics*, 107:1195–1225, Feb 2016.
- [738] Sheikh A Tahmid and Gennaro Notomista. Value iteration for learning concurrently executable robotic control tasks. *arXiv preprint arXiv:2504.01174*, 2025.
- [739] Cheng Tan, Dongxin Lyu, Siyuan Li, Zhangyang Gao, Jingxuan Wei, Siqi Ma, Zicheng Liu, and Stan Z Li. Peer review as a multi-turn and long-context dialogue with role-based interactions. *arXiv preprint arXiv:2406.05688*, 2024.
- [740] Hongming Tan, Shaoxiong Zhan, Fengwei Jia, Hai-Tao Zheng, and Wai Kin Chan. A hierarchical framework for measuring scientific paper innovation via large language models. *arXiv preprint arXiv:2504.14620*, 2025.
- [741] Xiangru Tang, Howard Dai, Elizabeth Knight, Fang Wu, Yunyang Li, Tianxiao Li, and Mark Gerstein. A survey of generative ai for de novo drug design: new frontiers in molecule and protein generation. *Briefings in Bioinformatics*, 25(4):bbae338, Jul 2024.
- [742] Xiangru Tang, Xingyao Zhang, Yanjun Shao, Jie Wu, Yilun Zhao, Arman Cohan, Ming Gong, Dongmei Zhang, and Mark Gerstein. Step-back profiling: Distilling user history for personalized scientific writing. *arXiv preprint arXiv:2406.14275*, 2024.
- [743] Zecheng Tang, Keyan Zhou, Juntao Li, Baibei Ji, Jianye Hou, and Min Zhang. L-citeeval: Do long-context models truly leverage context for responding? *arXiv preprint arXiv:2410.02115*, 2024.
- [744] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.

- [745] Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024.
- [746] Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhubapatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*, 2024.
- [747] Illustrae Team. How to create accurate scientific illustrations with ai in 2025. May 2025. URL <https://illustrae.co/blog/how-to-create-accurate-scientific-illustrations-ai>.
- [748] Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun Xiao, Chenzhuang Du, Chonghua Liao, et al. Kimi k1.5: Scaling reinforcement learning with llms. *arXiv preprint arXiv:2501.12599*, 2025.
- [749] NovelSeek Team, Bo Zhang, Shiyang Feng, Xiangchao Yan, Jiakang Yuan, Zhiyin Yu, Xiaohan He, Songtao Huang, Shaowei Hou, Zheng Nie, et al. Novelseek: When agent becomes the scientist—building closed-loop system from hypothesis to verification. *arXiv preprint arXiv:2505.16938*, 2025.
- [750] Qwen Team. Qwq: Reflect deeply on the boundaries of the unknown. <https://qwenlm.github.io/blog/qwq-32b-preview/>, Dec 2024. Accessed: 2024-12-16.
- [751] Qwen Team. Qwq-32b: Embracing the power of reinforcement learning, March 2025. URL <https://qwenlm.github.io/blog/qwq-32b/>.
- [752] Adrian Tedford. Helping editors find reviewers. <https://www.elsevier.com/connect/helping-editors-find-reviewers>, Sep 2015.
- [753] Mike Thelwall and Abdallah Yaghi. Evaluating the predictive capacity of chatgpt for academic peer review outcomes across multiple platforms. *Scientometrics*, pages 1–23, Mar 2025.
- [754] Tingzhong Tian, Shuya Li, Ziting Zhang, Lin Chen, Ziheng Zou, Dan Zhao, and Jianyang Zeng. Benchmarking compound activity prediction for real-world drug discovery applications. *Communications Chemistry*, 7(1):127, Jun 2024.
- [755] Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel. Domain randomization for transferring deep neural networks from simulation to the real world. In *2017 IEEE/RSJ international conference on intelligent robots and systems (IROS)*, pages 23–30. IEEE, 2017.
- [756] Gary Tom, Stefan P Schmid, Sterling G Baird, Yang Cao, Kourosh Darvish, Han Hao, Stanley Lo, Sergio Pablo-García, Ella M Rajaonson, Marta Skreta, et al. Self-driving laboratories for chemistry and materials science. *Chemical Reviews*, 124(16):9633–9732, Aug 2024.
- [757] Song Tong, Kai Mao, Zhen Huang, Yukun Zhao, and Kaiping Peng. Automating psychological hypothesis generation with ai: when large language models meet causal graph. *Humanities and Social Sciences Communications*, 11(1):1–14, Jul 2024.
- [758] Alexander Tornede, Difan Deng, Theresa Eimer, Joseph Giovanelli, Aditya Mohan, Tim Ruhkopf, Sarah Segel, Daphne Theodorakopoulos, Tanja Tornede, Henning Wachsmuth, et al. Automl in the age of large language models: Current challenges, future opportunities and risks. *arXiv preprint arXiv:2306.08107*, 2023.

- [759] Vladislav Trifonov, Iaroslav Kononov, Daniil Sherki, Oleg Svidchenko, Aleksei Shpilman, and Ekaterina Muravleva. Ai-powered platform for scientific discovery. In *AI4X 2025 International Conference*, Jul .
- [760] Stacey Truex, Nathalie Baracaldo, Ali Anwar, Thomas Steinke, Heiko Ludwig, Rui Zhang, and Yi Zhou. A hybrid approach to privacy-preserving federated learning. In *Proceedings of the 12th ACM workshop on artificial intelligence and security*, pages 1–11, Aug 2019.
- [761] Joseph Tu, Hilda Hadan, Derrick M Wang, Sabrina A Sgandurra, Reza Hadi Mogavi, and Lennart E Nacke. Augmenting the author: Exploring the potential of ai collaboration in academic writing. *arXiv preprint arXiv:2404.16071*, 2024.
- [762] Keith Tyser, Ben Segev, Gaston Longhitano, Xin-Yu Zhang, Zachary Meeks, Jason Lee, Uday Garg, Nicholas Belsten, Avi Shporer, Madeleine Udell, et al. Ai-driven review systems: evaluating llms in scalable and bias-aware academic reviews. *arXiv preprint arXiv:2408.10365*, 2024.
- [763] Shiekh Zia Uddin, Sachin Vaidya, Shrish Choudhary, Zhuo Chen, Raafat K Salib, Luke Huang, Dirk R Englund, and Marin Soljačić. Ai-driven robotics for free-space optics. *arXiv preprint arXiv:2505.17985*, 2025.
- [764] Pietro Vischia. Ai-assisted design of experiments at the frontiers of computation: methods and new perspectives. *arXiv preprint arXiv:2501.04448*, 2025.
- [765] Juraj Vladika and Florian Matthes. Comparing knowledge sources for open-domain scientific claim verification. *arXiv preprint arXiv:2402.02844*, 2024.
- [766] Juraj Vladika and Florian Matthes. Improving health question answering with reliable and time-aware evidence retrieval. *arXiv preprint arXiv:2404.08359*, 2024.
- [767] Luka Vrečar, Joe Wells, and Fairouz Kamareddine. Towards semantic markup of mathematical documents via user interaction. In *International Conference on Intelligent Computer Mathematics*, pages 223–240. Springer, Jul 2024.
- [768] David Wadden, Kyle Lo, Lucy Lu Wang, Arman Cohan, Iz Beltagy, and Hannaneh Hajishirzi. Multivers: Improving scientific claim verification with weak supervision and full-document context. *arXiv preprint arXiv:2112.01640*, 2021.
- [769] David Wadden, Kejian Shi, Jacob Morrison, Aakanksha Naik, Shruti Singh, Nitzan Barzilay, Kyle Lo, Tom Hope, Luca Soldaini, Zejiang Shen, et al. Sciriff: A resource to enhance language model instruction-following over scientific literature. In *Neurips 2024 Workshop Foundation Models for Science: Progress, Opportunities, and Challenges*, Dec 2024.
- [770] Byron C Wallace, Sayantan Saha, Frank Soboczinski, and Iain J Marshall. Generating (factual?) narrative summaries of rcts: Experiments with neural multi-document summarization. *AMIA Summits on Translational Science Proceedings*, 2021:605, May 2021.
- [771] Yuwei Wan, Yixuan Liu, Aswathy Ajith, Clara Grazian, Bram Hoex, Wenjie Zhang, Chunyu Kit, Tong Xie, and Ian Foster. Sciqag: A framework for auto-generated science question answering dataset with fine-grained evaluation. *arXiv preprint arXiv:2405.09939*, 2024.
- [772] Chengshi Wang, Yeon-Ju Kim, Aikaterini Vriza, Rohit Batra, Arun Baskaran, Naisong Shan, Nan Li, Pierre Darancet, Logan Ward, Yuzi Liu, et al. Autonomous platform for solution processing of electronic polymers. *Nature communications*, 16(1):1498, 2025.

- [773] Dingzirui Wang, Longxu Dou, and Wanxiang Che. A survey on table-and-text hybridqa: Concepts, methods, challenges and future directions. *arXiv preprint arXiv:2212.13465*, 2022.
- [774] Dingzirui Wang, Longxu Dou, Xuanliang Zhang, Qingfu Zhu, and Wanxiang Che. Improving demonstration diversity by human-free fusing for text-to-sql. *arXiv preprint arXiv:2402.10663*, 2024.
- [775] Haiming Wang, Huajian Xin, Chuanyang Zheng, Lin Li, Zhengying Liu, Qingxing Cao, Yinya Huang, Jing Xiong, Han Shi, Enze Xie, et al. Lego-prover: Neural theorem proving with growing libraries. *arXiv preprint arXiv:2310.00656*, 2023.
- [776] Haiming Wang, Ye Yuan, Zhengying Liu, Jianhao Shen, Yichun Yin, Jing Xiong, Enze Xie, Han Shi, Yujun Li, Lin Li, et al. Dt-solver: Automated theorem proving with dynamic-tree sampling guided by proof-level value function. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12632–12646, Jul 2023.
- [777] Haiming Wang, Huajian Xin, Zhengying Liu, Wenda Li, Yinya Huang, Jianqiao Lu, Zhicheng Yang, Jing Tang, Jian Yin, Zhenguo Li, et al. Proving theorems recursively. *arXiv preprint arXiv:2405.14414*, 2024.
- [778] Haining Wang. A content-based novelty measure for scholarly publications: A proof of concept. In *International Conference on Information*, pages 409–420. Springer, Apr 2024.
- [779] Hanchen Wang, Tianfan Fu, Yuanqi Du, Wenhao Gao, Kexin Huang, Ziming Liu, Payal Chandak, Shengchao Liu, Peter Van Katwyk, Andreea Deac, et al. Scientific discovery in the age of artificial intelligence. *Nature*, 620(7972):47–60, Aug 2023.
- [780] Hanchen Wang, Yichun He, Paula P Coelho, Matthew Bucci, Abbas Nazir, Bob Chen, Linh Trinh, Serena Zhang, Kexin Huang, Vineethkrishna Chandrasekar, et al. Spatialagent: An autonomous ai agent for spatial biology. *bioRxiv*, pages 2025–04, Apr 2025.
- [781] Haorui Wang, Marta Skreta, Cher-Tian Ser, Wenhao Gao, Lingkai Kong, Felix Strieth-Kalthoff, Chenru Duan, Yuchen Zhuang, Yue Yu, Yanqiao Zhu, et al. Efficient evolutionary search over chemical space with large language models. *arXiv preprint arXiv:2406.16976*, 2024.
- [782] Haoyu Wang, Yujia Fu, Zhu Zhang, Shuo Wang, Zirui Ren, Xiaorong Wang, Zhili Li, Chaoqun He, Bo An, Zhiyuan Liu, et al. Llm \times mapreduce-v2: Entropy-driven convolutional test-time scaling for generating long-form articles from extremely long resources. *arXiv preprint arXiv:2504.05732*, 2025.
- [783] Izia Xiaoxiao Wang, Xihan Wu, Edith Coates, Min Zeng, Jiexin Kuang, Siliang Liu, Mengyang Qiu, and Jungyeul Park. Neural automated writing evaluation with corrective feedback. *arXiv preprint arXiv:2402.17613*, 2024.
- [784] Jiashuo Wang, Yang Xiao, Yanran Li, Changhe Song, Chunpu Xu, Chenhao Tan, and Wenjie Li. Towards a client-centered assessment of llm therapists by client simulation. *arXiv preprint arXiv:2406.12266*, 2024.
- [785] Lei Wang, Yi Hu, Jiabang He, Xing Xu, Ning Liu, Hui Liu, and Heng Tao Shen. T-sciq: Teaching multimodal chain-of-thought reasoning via mixed large language model signals for science question answering. *arXiv preprint arXiv:2305.03453*, 2023.
- [786] Linghe Wang, Minhwa Lee, Ross Volkov, Luan Tuyen Chau, and Dongyeop Kang. Scholawrite: A dataset of end-to-end scholarly writing process. *arXiv preprint arXiv:2502.02904*, 2025.

- [787] Lucy Lu Wang, Jay DeYoung, and Byron Wallace. Overview of mslr2022: A shared task on multi-document summarization for literature reviews. In *Proceedings of the third workshop on scholarly document processing*, Oct 2022.
- [788] Pancheng Wang, Shasha Li, Haifang Zhou, Jintao Tang, and Ting Wang. Toc-rwg: Explore the combination of topic model and citation information for automatic related work generation. *Ieee Access*, 8:13043–13055, Dec 2019.
- [789] Pancheng Wang, Shasha Li, Kunyuan Pang, Liangliang He, Dong Li, Jintao Tang, and Ting Wang. Multi-document scientific summarization from a knowledge graph-centric view. *arXiv preprint arXiv:2209.04319*, 2022.
- [790] Pancheng Wang, Shasha Li, Dong Li, Kehan Long, Jintao Tang, and Ting Wang. Disentangling instructive information from ranked multiple candidates for multi-document scientific summarization. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 2028–2037, Jul 2024.
- [791] Peng Wang, Yongheng Zhang, Hao Fei, Qiguang Chen, Yukai Wang, Jiasheng Si, Wenpeng Lu, Min Li, and Libo Qin. S3 agent: Unlocking the power of vllm for zero-shot multi-modal sarcasm detection. *ACM Transactions on Multimedia Computing, Communications and Applications*, Aug 2024.
- [792] Peng Wang, Ruihan Tao, Qiguang Chen, Mengkang Hu, and Libo Qin. X-webagentbench: A multilingual interactive web benchmark for evaluating global agentic system. *arXiv preprint arXiv:2505.15372*, 2025.
- [793] Qingyun Wang, Doug Downey, Heng Ji, and Tom Hope. Scimon: Scientific inspiration machines optimized for novelty. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 279–299, Aug 2024.
- [794] Runxiang Wang, Boxiao Wang, Kai Li, Yifan Zhang, and Jian Cheng. Drsr: Llm based scientific equation discovery with dual reasoning from data and experience. *arXiv preprint arXiv:2506.04282*, 2025.
- [795] Siyuan Wang, James R Foulds, Md Osman Gani, and Shimei Pan. Llm-based corroborating and refuting evidence retrieval for scientific claim verification. *arXiv preprint arXiv:2503.07937*, 2025.
- [796] Wenxiao Wang, Lihui Gu, Liye Zhang, Yunxiang Luo, Yi Dai, Chen Shen, Liang Xie, Binbin Lin, Xiaofei He, and Jieping Ye. Scipip: An llm-based scientific paper idea proposer. *arXiv preprint arXiv:2410.23166*, 2024.
- [797] Wenxuan Wang, Zizhan Ma, Zheng Wang, Chenghan Wu, Wenting Chen, Xiang Li, and Yixuan Yuan. A survey of llm-based agents in medicine: How far are we from baymax? *arXiv preprint arXiv:2502.11211*, 2025.
- [798] Xiaoxuan Wang, Ziniu Hu, Pan Lu, Yanqiao Zhu, Jieyu Zhang, Satyen Subramaniam, Arjun R Loomba, Shichang Zhang, Yizhou Sun, and Wei Wang. Scibench: Evaluating college-level scientific problem-solving abilities of large language models. *arXiv preprint arXiv:2307.10635*, 2023.
- [799] Xingyao Wang, Boxuan Li, Yufan Song, Frank F Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan, Yueqi Song, Bowen Li, Jaskirat Singh, et al. Opendevin: An open platform for ai software developers as generalist agents. *arXiv preprint arXiv:2407.16741*, 2024.

- [800] Xingyao Wang, Boxuan Li, Yufan Song, Frank F Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan, Yueqi Song, Bowen Li, Jaskirat Singh, et al. Openhands: An open platform for ai software developers as generalist agents. *arXiv preprint arXiv:2407.16741*, 2024.
- [801] Yao Wang, Mingxuan Cui, and Arthur Jiang. Enabling ai scientists to recognize innovation: A domain-agnostic algorithm for assessing novelty. *arXiv preprint arXiv:2503.01508*, 2025.
- [802] Yidong Wang, Qi Guo, Wenjin Yao, Hongbo Zhang, Xin Zhang, Zhen Wu, Meishan Zhang, Xinyu Dai, Qingsong Wen, Wei Ye, et al. Autosurvey: Large language models can automatically write surveys. *Advances in Neural Information Processing Systems*, 37:115119–115145, Dec 2024.
- [803] Yifan Wang, David Stevens, Pranay Shah, Wenwen Jiang, Miao Liu, Xu Chen, Robert Kuo, Na Li, Boying Gong, Daniel Lee, et al. Model-in-the-loop (milo): Accelerating multimodal ai data annotation with llms. *arXiv preprint arXiv:2409.10702*, 2024.
- [804] Yixuan Wang, Baoxin Wang, Yijun Liu, Dayong Wu, and Wanxiang Che. Lm-combiner: A contextual rewriting model for chinese grammatical error correction. *arXiv preprint arXiv:2403.17413*, 2024.
- [805] Yixuan Wang, Baoxin Wang, Yijun Liu, Qingfu Zhu, Dayong Wu, and Wanxiang Che. Improving grammatical error correction via contextual data augmentation. *arXiv preprint arXiv:2406.17456*, 2024.
- [806] Yizheng Wang, Jinshuai Bai, Zhongya Lin, Qimin Wang, Cosmin Animescu, Jia Sun, Mohammad Sadegh Eshaghi, Yuantong Gu, Xi-Qiao Feng, Xiaoying Zhuang, et al. Artificial intelligence for partial differential equations in computational mechanics: A review. *arXiv preprint arXiv:2410.19843*, 2024.
- [807] Yongzhen Wang, Xiaozhong Liu, and Zheng Gao. Neural related work summarization with a joint context-driven attention mechanism. *arXiv preprint arXiv:1901.09492*, 2019.
- [808] Yubo Wang, Xueguang Ma, Ping Nie, Huaye Zeng, Zhiheng Lyu, Yuxuan Zhang, Benjamin Schneider, Yi Lu, Xiang Yue, and Wenhui Chen. Scholarcopilot: Training large language models for academic writing with accurate citations. *arXiv preprint arXiv:2504.00824*, 2025.
- [809] Zifeng Wang, Benjamin Danek, and Jimeng Sun. Bidsa-1k: Benchmarking data science agents for biomedical research. *arXiv preprint arXiv:2505.16100*, 2025.
- [810] Zilong Wang, Hao Zhang, Chun-Liang Li, Julian Martin Eisenschlos, Vincent Perot, Zifeng Wang, Lesly Miculicich, Yasuhisa Fujii, Jingbo Shang, Chen-Yu Lee, et al. Chain-of-table: Evolving tables in the reasoning chain for table understanding. In *The Twelfth International Conference on Learning Representations*, Jan 2024. URL <https://openreview.net/forum?id=4L0xnS4GQM>.
- [811] Ziqi Wang, Le Hou, Tianjian Lu, Yuexin Wu, Yunxuan Li, Hongkun Yu, and Heng Ji. Enabling language models to implicitly learn self-improvement. *arXiv preprint arXiv:2310.00898*, 2023.
- [812] Zirui Wang, Mengzhou Xia, Luxi He, Howard Chen, Yitao Liu, Richard Zhu, Kaiqu Liang, Xindi Wu, Haotian Liu, Sadhika Malladi, Alexis Chevalier, Sanjeev Arora, and Danqi Chen. Charxiv: Charting gaps in realistic chart understanding in multimodal llms. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang, editors, *Advances in Neural Information Processing Systems*, volume 37, pages 113569–113697. Curran Associates, Inc., Dec 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/cdf6f8e9fd9aeaf79b6024caec24f15b-Paper-Datasets_and_Benchmarks_Track.pdf.

- [813] William Watson and Lawrence Yong. Directed criteria citation recommendation and ranking through link prediction. *arXiv preprint arXiv:2403.18855*, 2024.
- [814] Jane Webster and Richard T Watson. Analyzing the past to prepare for the future: Writing a literature review. *MIS quarterly*, pages xiii–xxiii, Jun 2002.
- [815] Jason Wei, Zhiqing Sun, Spencer Papay, Scott McKinney, Jeffrey Han, Isa Fulford, Hyung Won Chung, Alex Tachard Passos, William Fedus, and Amelia Glaese. Browsecomp: A simple yet challenging benchmark for browsing agents. *arXiv preprint arXiv:2504.12516*, 2025.
- [816] Haomin Wen, Zhenjie Wei, Yan Lin, Jiyuan Wang, Yuxuan Liang, and Huaiyu Wan. Overleafcopilot: Empowering academic writing in overleaf with large language models. *arXiv preprint arXiv:2403.09733*, 2024.
- [817] Jiaxin Wen, Chenglei Si, Yueh-han Chen, He He, and Shi Feng. Predicting empirical ai research outcomes with language models. *arXiv preprint arXiv:2506.00794*, 2025.
- [818] Ju Wen and Lan Yi. Are plain language summaries more readable than scientific abstracts? evidence from six biomedical and life sciences journals. *Public Understanding of Science*, 34(1):114–126, May 2025.
- [819] Yingting Wen and Sandra Laporte. Experiential narratives in marketing: A comparison of generative ai and human content. *Journal of Public Policy & Marketing*, 44(3):392–410, Oct 2025. doi: 10.1177/07439156241297973. URL <https://doi.org/10.1177/07439156241297973>.
- [820] Yixuan Weng, Minjun Zhu, Guangsheng Bao, Hongbo Zhang, Jindong Wang, Yue Zhang, and Linyi Yang. Cyclereviewer: Improving automated research via automated review. *arXiv preprint arXiv:2411.00816*, 2024.
- [821] Nicole E Wheeler. Responsible ai in biotechnology: balancing discovery, innovation and biosecurity risks. *Frontiers in Bioengineering and Biotechnology*, 13:1537471, Feb 2025.
- [822] Dustin Wright and Isabelle Augenstein. Citeworth: Cite-worthiness detection for improved scientific document understanding. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 1796–1807, May 2021.
- [823] Chaoyi Wu, Jiayu Lei, Qiaoyu Zheng, Weike Zhao, Weixiong Lin, Xiaoman Zhang, Xiao Zhou, Ziheng Zhao, Ya Zhang, Yanfeng Wang, et al. Can gpt-4v (ision) serve medical applications? case studies on gpt-4v for multimodal medical diagnosis. *arXiv preprint arXiv:2310.09909*, Oct 2023.
- [824] Jialong Wu, Baixuan Li, Runnan Fang, Wenbiao Yin, Liwen Zhang, Zhengwei Tao, Dingchu Zhang, Zekun Xi, Yong Jiang, Pengjun Xie, et al. Webdancer: Towards autonomous information seeking agency. *arXiv preprint arXiv:2505.22648*, 2025.
- [825] Jian Wu, Jiayu Zhang, Dongyuan Li, Linyi Yang, Aoxiao Zhong, Renhe Jiang, Qingsong Wen, and Yue Zhang. Lag: Llm agents for leaderboard auto generation on demanding. *arXiv preprint arXiv:2502.18209*, 2025.
- [826] Jinxuan Wu, Wenhan Chao, Xian Zhou, and Zhunchen Luo. Characterizing and verifying scientific claims: Qualitative causal structure is all you need. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 13428–13439, Dec 2023.

- [827] Junfeng Wu, Jing He, Hai Liu, Zhaoqi Zheng, Yichen Cao, Xingguo Chen, Bingjie Zou, Ruiping Zou, Guohua Zhou, David Sturgess, et al. From literature to lab: Hardware-independent autonomous chemical synthesis with reinforcement learning. In *Companion Proceedings of the ACM on Web Conference 2025*, pages 2923–2926, May 2025.
- [828] Shican Wu, Xiao Ma, Dehui Luo, Lulu Li, Xiangcheng Shi, Xin Chang, Xiaoyun Lin, Ran Luo, Chunlei Pei, Changying Du, et al. Automated review generation method based on large language models. *arXiv preprint arXiv:2407.20906*, 2024.
- [829] Siyang Wu, Honglin Bao, Nadav Kunievsky, and James A Evans. Introspective growth: Automatically advancing llm expertise in technology judgment. *arXiv preprint arXiv:2505.12452*, 2025.
- [830] Siyuan Wu, Leong Hou U, Sourav S Bhowmick, and Wolfgang Gatterbauer. Pistis: A conflict of interest declaration and detection system for peer review management. In *Proceedings of the 2018 International Conference on Management of Data*, pages 1713–1716, May 2018.
- [831] Wenqing Wu, Chengzhi Zhang, Tong Bao, and Yi Zhao. Sc4anm: Identifying optimal section combinations for automated novelty prediction in academic papers. *Expert Systems with Applications*, page 126778, May 2025.
- [832] Xianjie Wu, Jian Yang, Linzheng Chai, Ge Zhang, Jiaheng Liu, Xeron Du, Di Liang, Daixin Shu, Xianfu Cheng, Tianzhen Sun, et al. Tablebench: A comprehensive and complex benchmark for table question answering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 25497–25506, Apr 2025.
- [833] Yushuai Wu, Ting Zhang, Hao Zhou, Hainan Wu, Hanwen Sunchu, Lei Hu, Xiaofang Chen, Suyuan Zhao, Gaochao Liu, Chao Sun, et al. Deepcre: Transforming drug r&d via ai-driven cross-drug response evaluation. *arXiv preprint arXiv:2403.03768*, 2024.
- [834] Amelie Wühl, Yarik Menchaca Resendiz, Lara Grimminger, and Roman Klinger. What makes medical claims (un) verifiable? analyzing entity and relation properties for fact verification. *arXiv preprint arXiv:2402.01360*, 2024.
- [835] Amelie Wühl, Dustin Wright, Roman Klinger, and Isabelle Augenstein. Understanding fine-grained distortions in reports of scientific findings. *arXiv preprint arXiv:2402.12431*, 2024.
- [836] X.AI. Grok-2 beta release. <https://x.ai/blog/grok-2>, Dec 2024. Accessed: 2024-12-16.
- [837] Renqiu Xia, Bo Zhang, Hancheng Ye, Xiangchao Yan, Qi Liu, Hongbin Zhou, Zijun Chen, Peng Ye, Min Dou, Botian Shi, et al. Chartx & chartvlm: A versatile benchmark and foundation model for complicated chart reasoning. *arXiv preprint arXiv:2402.12185*, 2024.
- [838] Yutong Xia, Ao Qu, Yunhan Zheng, Yihong Tang, Dingyi Zhuang, Yuxuan Liang, Shenhao Wang, Cathy Wu, Lijun Sun, Roger Zimmermann, et al. Reimagining urban science: Scaling causal inference with large language models. *arXiv preprint arXiv:2504.12345*, 2025.
- [839] Yanzheng Xiang, Hanqi Yan, Shuyin Ouyang, Lin Gui, and Yulan He. Scireplicate-bench: Benchmarking llms in agent-driven algorithmic reproduction from research papers. *arXiv preprint arXiv:2504.00255*, 2025.

- [840] Yihang Xiao, Jinyi Liu, Yan Zheng, Xiaohan Xie, Jianye Hao, Mingzhi Li, Ruitao Wang, Fei Ni, Yuxiao Li, Jintian Luo, et al. Cellagent: An llm-driven multi-agent framework for automated single-cell data analysis. *arXiv preprint arXiv:2407.09811*, 2024.
- [841] Huajian Xin, Daya Guo, Zhihong Shao, Zhizhou Ren, Qihao Zhu, Bo Liu, Chong Ruan, Wenda Li, and Xiaodan Liang. Deepseek-prover: Advancing theorem proving in llms through large-scale synthetic data. *arXiv preprint arXiv:2405.14333*, 2024.
- [842] Junjie Xing, Yeye He, Mengyu Zhou, Haoyu Dong, Shi Han, Dongmei Zhang, and Surajit Chaudhuri. Table-llm-specialist: Language model specialists for tables using iterative generator-validator fine-tuning. *arXiv preprint arXiv:2410.12164*, 2024.
- [843] Guangzhi Xiong, Eric Xie, Corey Williams, Myles Kim, Amir Hassan Shariatmadari, Sikun Guo, Stefan Bekiranov, and Aidong Zhang. Toward reliable biomedical hypothesis generation: Evaluating truthfulness and hallucination in large language models. *arXiv preprint arXiv:2505.14599*, 2025.
- [844] Ruibin Xiong, Yimeng Chen, Dmitrii Khizbullin, Mingchen Zhuge, and Jürgen Schmidhuber. Beyond outlining: Heterogeneous recursive planning for adaptive long-form writing with language models. *arXiv preprint arXiv:2503.08275*, 2025.
- [845] Wanghan Xu, Xiangyu Zhao, Yuhao Zhou, Xiaoyu Yue, Ben Fei, Fenghua Ling, Wenlong Zhang, and Lei Bai. Earthse: A benchmark evaluating earth scientific exploration capability for large language models. *arXiv preprint arXiv:2505.17139*, 2025.
- [846] Ziyang Xu. Patterns and purposes: A cross-journal analysis of ai tool usage in academic writing. *arXiv preprint arXiv:2502.00632*, 2025.
- [847] Kevin G Yager. Towards a science exocortex. *Digital Discovery*, 3(10):1933–1957, Aug 2024.
- [848] Yutaro Yamada, Robert Tjarko Lange, Cong Lu, Shengran Hu, Chris Lu, Jakob Foerster, Jeff Clune, and David Ha. The ai scientist-v2: Workshop-level automated scientific discovery via agentic tree search. *arXiv preprint arXiv:2504.08066*, 2025.
- [849] Xiangchao Yan, Shiyang Feng, Jiakang Yuan, Renqiu Xia, Bin Wang, Bo Zhang, and Lei Bai. Surveyforge: On the outline heuristics, memory-driven generation, and multi-dimensional evaluation for automated survey writing. *arXiv preprint arXiv:2503.04629*, 2025.
- [850] Yibo Yan, Shen Wang, Jiahao Huo, Jingheng Ye, Zhendong Chu, Xuming Hu, Philip S Yu, Carla Gomes, Bart Selman, and Qingsong Wen. Position: Multimodal large language models can significantly advance scientific reasoning. *arXiv preprint arXiv:2502.02871*, 2025.
- [851] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- [852] Hsin-Jung Yang, Mahsa Khosravi, Benjamin Walt, Girish Krishnan, and Soumik Sarkar. Zero-shot sim-to-real transfer for reinforcement learning-based visual servoing of soft continuum arms. *arXiv preprint arXiv:2504.16916*, 2025.

- [853] Jianke Yang, Manu Bhat, Bryan Hu, Yadi Cao, Nima Dehmamy, Robin Walters, and Rose Yu. Discovering symbolic differential equations with symmetry invariants. *arXiv preprint arXiv:2505.12083*, 2025.
- [854] John Jeongseok Yang and Sang-Hyun Hwang. Transforming hematological research documentation with large language models: an approach to scientific writing and data analysis. *Blood research*, 60(1):1–11, Mar 2025.
- [855] Lin Yang, Shawn Xu, Andrew Sellergren, Timo Kohlberger, Yuchen Zhou, Ira Ktena, Atilla Kiraly, Faruk Ahmed, Farhad Hormozdiari, Tiam Jaroensri, et al. Advancing multimodal medical capabilities of gemini. *arXiv preprint arXiv:2405.03162*, 2024.
- [856] Rui Yang, Boming Yang, Aosong Feng, Sixun Ouyang, Moritz Blum, Tianwei She, Yuang Jiang, Freddy Lecue, Jinghui Lu, and Irene Li. Graphusion: a rag framework for knowledge graph construction with a global perspective. *arXiv preprint arXiv:2410.17600*, 2024.
- [857] Yifei Yang, Runhan Shi, Zuchao Li, Shu Jiang, Yang Yang, Bao-Liang Lu, and Hai Zhao. Batgpt-chem: A foundation large model for chemical engineering. Apr 2024.
- [858] Yue Yang, MingKang Chen, Qihua Liu, Mengkang Hu, Qiguang Chen, Gengrui Zhang, Shuyue Hu, Guangtao Zhai, Yu Qiao, Yu Wang, et al. Truly assessing fluid intelligence of large language models through dynamic reasoning evaluation. *arXiv preprint arXiv:2506.02648*, 2025.
- [859] Zhaorui Yang, Bo Pan, Han Wang, Yiyao Wang, Xingyu Liu, Minfeng Zhu, Bo Zhang, and Wei Chen. Multimodal deepresearcher: Generating text-chart interleaved reports from scratch with agentic framework. *arXiv preprint arXiv:2506.02454*, 2025.
- [860] 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.
- [861] Zonglin Yang, Wanhao Liu, Ben Gao, Tong Xie, Yuqiang Li, Wanli Ouyang, Soujanya Poria, Erik Cambria, and Dongzhan Zhou. MOOSE-chem: Large language models for rediscovering unseen chemistry scientific hypotheses. In *The Thirteenth International Conference on Learning Representations*, Jan 2025. URL <https://openreview.net/forum?id=X9OfMNNepI>.
- [862] Zukang Yang, Zixuan Zhu, and Jennifer Zhu. Curiousllm: Elevating multi-document question answering with llm-enhanced knowledge graph reasoning. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 3: Industry Track)*, pages 274–286, Jan 2025.
- [863] Nicolas Yax, Hernán Anlló, and Stefano Palminteri. Studying and improving reasoning in humans and machines. *Communications Psychology*, 2(1):51, Jun 2024.
- [864] Mang Ye, Xiuwen Fang, Bo Du, Pong C Yuen, and Dacheng Tao. Heterogeneous federated learning: State-of-the-art and research challenges. *ACM Computing Surveys*, 56(3):1–44, Oct 2023.
- [865] Rui Ye, Xianghe Pang, Jingyi Chai, Jiaao Chen, Zhenfei Yin, Zhen Xiang, Xiaowen Dong, Jing Shao, and Siheng Chen. Are we there yet? revealing the risks of utilizing large language models in scholarly peer review. *arXiv preprint arXiv:2412.01708*, 2024.

- [866] Yanpeng Ye, Jie Ren, Shaozhou Wang, Yuwei Wan, Imran Razzak, Bram Hoex, Haofen Wang, Tong Xie, and Wenjie Zhang. Construction and application of materials knowledge graph in multidisciplinary materials science via large language model. *Advances in Neural Information Processing Systems*, 37: 56878–56897, Dec 2024.
- [867] Ho Yin, Ting-Yao Hsu, Jiyou Min, Sungchul Kim, Ryan A Rossi, Tong Yu, Hyunggu Jung, Ting-Hao‘Kenneth’ Huang, et al. Understanding how paper writers use ai-generated captions in figure caption writing. *arXiv preprint arXiv:2501.06317*, 2025.
- [868] Ming Yin, Yuanhao Qu, Dyllan Liu, Ling Yang, Le Cong, and Mengdi Wang. Genome-bench: A scientific reasoning benchmark from real-world expert discussions. *bioRxiv*, pages 2025–06, 2025.
- [869] Chengzhang Yu, Yiming Zhang, Zhixin Liu, Zenghui Ding, Yining Sun, and Zhanpeng Jin. Frame: Feedback-refined agent methodology for enhancing medical research insights. *arXiv preprint arXiv:2505.04649*, 2025.
- [870] Fei Yu, Hongbo Zhang, Prayag Tiwari, and Benyou Wang. Natural language reasoning, a survey. *ACM Comput. Surv.*, 56(12), October 2024. ISSN 0360-0300. doi: 10.1145/3664194. URL <https://doi.org/10.1145/3664194>.
- [871] Haofei Yu, Zhaochen Hong, Zirui Cheng, Kunlun Zhu, Keyang Xuan, Jinwei Yao, Tao Feng, and Jiaxuan You. Researchtown: Simulator of human research community. *arXiv preprint arXiv:2412.17767*, 2024.
- [872] Hengjie Yu and Yaochu Jin. Unlocking the potential of ai researchers in scientific discovery: What is missing? *arXiv preprint arXiv:2503.05822*, 2025.
- [873] Jianxiang Yu, Jiaqi Tan, Zichen Ding, Jiapeng Zhu, Jiahao Li, Yao Cheng, Qier Cui, Yunshi Lan, and Xiang Li. Seagraph: Unveiling the whole story of paper review comments. *arXiv preprint arXiv:2412.11939*, 2024.
- [874] Jing Yu, Yuqi Tang, Kehua Feng, Mingyang Rao, Lei Liang, Zhiqiang Zhang, Mengshu Sun, Wen Zhang, Qiang Zhang, Keyan Ding, et al. Scicueval: A comprehensive dataset for evaluating scientific context understanding in large language models. *arXiv preprint arXiv:2505.15094*, 2025.
- [875] Luyao Yu, Qi Zhang, Chongyang Shi, An Lao, and Liang Xiao. Reinforced subject-aware graph neural network for related work generation. In *International Conference on Knowledge Science, Engineering and Management*, pages 201–213. Springer, Jul 2024.
- [876] Sungduk Yu, Man Luo, Avinash Madasu, Vasudev Lal, and Phillip Howard. Is your paper being reviewed by an llm? investigating ai text detectability in peer review. *arXiv preprint arXiv:2410.03019*, 2024.
- [877] Jiakang Yuan, Xiangchao Yan, Botian Shi, Tao Chen, Wanli Ouyang, Bo Zhang, Lei Bai, Yu Qiao, and Bowen Zhou. Dolphin: Closed-loop open-ended auto-research through thinking, practice, and feedback. *arXiv preprint arXiv:2501.03916*, 2025.
- [878] Shuzhou Yuan and Michael Färber. Hallucinations can improve large language models in drug discovery. *arXiv preprint arXiv:2501.13824*, 2025.
- [879] Weizhe Yuan and Pengfei Liu. Kid-review: knowledge-guided scientific review generation with oracle pre-training. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 11639–11647, Jun 2022.

- [880] Weizhe Yuan, Pengfei Liu, and Graham Neubig. Can we automate scientific reviewing? *Journal of Artificial Intelligence Research*, 75:171–212, Sep 2022.
- [881] Wenhao Yuan, Guangyao Chen, Zhilong Wang, and Fengqi You. Empowering generalist material intelligence with large language models. *Advanced Materials*, page 2502771, May 2025.
- [882] Miguel Zabaleta and Joel Lehman. Simulating tabular datasets through llms to rapidly explore hypotheses about real-world entities. *arXiv preprint arXiv:2411.18071*, 2024.
- [883] César Zamudio, Jamie L Grigsby, and Meg Michelsen. Raise: A new method to develop experimental stimuli for advertising research with image generative artificial intelligence. *Journal of Advertising*, pages 1–16, Jan 2024.
- [884] Fengzhu Zeng and Wei Gao. Prompt to be consistent is better than self-consistent? few-shot and zero-shot fact verification with pre-trained language models. *arXiv preprint arXiv:2306.02569*, 2023.
- [885] Qi Zeng, Mankeerat Sidhu, Hou Pong Chan, Lu Wang, and Heng Ji. Meta-review generation with checklist-guided iterative introspection. *arXiv preprint arXiv:2305.14647*, 2023.
- [886] Qi Zeng, Mankeerat Sidhu, Ansel Blume, Hou Pong Chan, Lu Wang, and Heng Ji. Scientific opinion summarization: Paper meta-review generation dataset, methods, and evaluation. In *Artificial Intelligence for Research and Democracy: First International Workshop, AI4Research 2024, and 4th International Workshop, DemocrAI 2024, Held in Conjunction with IJCAI 2024, Jeju, South Korea, August 5, 2024, Proceedings*, page 20. Springer Nature, Jun 2024.
- [887] Yongchao Zeng, Calum Brown, Joanna Raymond, Mohamed Byari, Ronja Hotz, and Mark Rounsevell. Exploring the opportunities and challenges of using large language models to represent institutional agency in land system modelling. *EGUsphere*, 2024:1–35, Mar 2024.
- [888] Chen Zhang, Yu Xie, Hang Bai, Bin Yu, Weihong Li, and Yuan Gao. A survey on federated learning. *Knowledge-Based Systems*, 216:106775, Mar 2021.
- [889] Dan Zhang, Ziniu Hu, Sining Zhoubian, Zhengxiao Du, Kaiyu Yang, Zihan Wang, Yisong Yue, Yuxiao Dong, and Jie Tang. Sciglm: Training scientific language models with self-reflective instruction annotation and tuning. *arXiv preprint arXiv:2401.07950*, 2024.
- [890] Daoze Zhang, Zhijian Bao, Sihang Du, Zhiyi Zhao, Kuangling Zhang, Dezheng Bao, and Yang Yang. Re²: A consistency-ensured dataset for full-stage peer review and multi-turn rebuttal discussions. *arXiv preprint arXiv:2505.07920*, 2025.
- [891] Fanjin Zhang, Shijie Shi, Yifan Zhu, Bo Chen, Yukuo Cen, Jifan Yu, Yelin Chen, Lulu Wang, Qingfei Zhao, Yuqing Cheng, Tianyi Han, Yuwei An, Dan Zhang, Weng Lam Tam, Kun Cao, Yunhe Pang, Xinyu Guan, Huihui Yuan, Jian Song, Xiaoyan Li, Yuxiao Dong, and Jie Tang. Oag-bench: A human-curated benchmark for academic graph mining. *arXiv preprint arXiv:2402.15810*, 2024.
- [892] Jintian Zhang, Xin Xu, Ningyu Zhang, Ruibo Liu, Bryan Hooi, and Shumin Deng. Exploring collaboration mechanisms for llm agents: A social psychology view. *arXiv preprint arXiv:2310.02124*, 2023.
- [893] Lei Zhang, Yuge Zhang, Kan Ren, Dongsheng Li, and Yuqing Yang. Mlcpilot: Unleashing the power of large language models in solving machine learning tasks. *arXiv preprint arXiv:2304.14979*, 2023.

- [894] Leixin Zhang, Steffen Eger, Yinjie Cheng, Weihe Zhai, Jonas Belouadi, Christoph Leiter, Simone Paolo Ponzetto, Fahimeh Moafian, and Zhixue Zhao. Scimage: How good are multimodal large language models at scientific text-to-image generation? *arXiv preprint arXiv:2412.02368*, 2024.
- [895] Michael JQ Zhang and Eunsol Choi. Clarify when necessary: Resolving ambiguity through interaction with lms. *arXiv preprint arXiv:2311.09469*, 2023.
- [896] Ming Zhang, Yujiong Shen, Zelin Li, Huayu Sha, Binze Hu, Yuhui Wang, Chenhao Huang, Shichun Liu, Jingqi Tong, Changhao Jiang, et al. Llmeval-med: A real-world clinical benchmark for medical llms with physician validation. *arXiv preprint arXiv:2506.04078*, 2025.
- [897] Shujian Zhang, Chengyue Gong, Lemeng Wu, Xingchao Liu, and Mingyuan Zhou. Automl-gpt: Automatic machine learning with gpt. *arXiv preprint arXiv:2305.02499*, 2023.
- [898] Tao Zhang, Zhenhai Liu, Yong Xin, and Yongjun Jiao. Mooseagent: A llm based multi-agent framework for automating moose simulation. *arXiv preprint arXiv:2504.08621*, 2025.
- [899] Tianmai M Zhang and Neil F Abernethy. Reviewing scientific papers for critical problems with reasoning llms: Baseline approaches and automatic evaluation. *arXiv preprint arXiv:2505.23824*, 2025.
- [900] Xiaocheng Zhang, Xi Wang, Yifei Lu, Zhuangzhuang Ye, Jianing Wang, Mengjiao Bao, Peng Yan, and Xiaohong Su. Augmenting the veracity and explanations of complex fact checking via iterative self-revision with llms. *arXiv preprint arXiv:2410.15135*, 2024.
- [901] Xingjian Zhang, Yutong Xie, Jin Huang, Jinge Ma, Zhaoying Pan, Qijia Liu, Ziyang Xiong, Tolga Ergen, Dongsub Shim, Honglak Lee, et al. Massw: A new dataset and benchmark tasks for ai-assisted scientific workflows. *arXiv preprint arXiv:2406.06357*, 2024.
- [902] Xinyu Zhang, Zongming Ni, Billy Fanady, Feibei Chen, Zijian Chen, Zixuan Wang, Guofei Chen, Zhengda He, Yang Bai, and Haitao Zhao. Chatgpt-assisted rational design for iterative performance optimization of perovskite solar cells. *Available at SSRN 5127472*, Feb .
- [903] Xuan Zhang and Wei Gao. Towards llm-based fact verification on news claims with a hierarchical step-by-step prompting method. *arXiv preprint arXiv:2310.00305*, 2023.
- [904] Xuanliang Zhang, Dingzirui Wang, Longxu Dou, Qingfu Zhu, and Wanxiang Che. A survey of table reasoning with large language models. *Frontiers of Computer Science*, 19(9):199348, Jan 2025.
- [905] Yanbo Zhang, Sumeer A Khan, Adnan Mahmud, Huck Yang, Alexander Lavin, Michael Levin, Jeremy Frey, Jared Dunnmon, James Evans, Alan Bundy, et al. Advancing the scientific method with large language models: From hypothesis to discovery. *arXiv preprint arXiv:2505.16477*, 2025.
- [906] Yang Zhang, Yufei Wang, Kai Wang, Quan Z Sheng, Lina Yao, Adnan Mahmood, Wei Emma Zhang, and Rongying Zhao. When large language models meet citation: A survey. *arXiv preprint arXiv:2309.09727*, 2023.
- [907] Yongheng Zhang, Qiguang Chen, Min Li, Wanxiang Che, and Libo Qin. Autocap: Towards automatic cross-lingual alignment planning for zero-shot chain-of-thought. *arXiv preprint arXiv:2406.13940*, 2024.

- [908] Yongheng Zhang, Qiguang Chen, Jingxuan Zhou, Peng Wang, Jiasheng Si, Jin Wang, Wenpeng Lu, and Libo Qin. Wrong-of-thought: An integrated reasoning framework with multi-perspective verification and wrong information. *arXiv preprint arXiv:2410.04463*, 2024.
- [909] Yongheng Zhang, Xu Liu, Ruoxi Zhou, Qiguang Chen, Hao Fei, Wenpeng Lu, and Libo Qin. Cchallenge: A novel benchmark for joint cross-lingual and cross-modal hallucinations detection in large language models. *arXiv preprint arXiv:2505.19108*, 2025.
- [910] 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*, 2024.
- [911] Yunxiang Zhang, Muhammad Khalifa, Shitanshu Bhushan, Grant D Murphy, Lajanugen Logeswaran, Jaekyeom Kim, Moontae Lee, Honglak Lee, and Lu Wang. Mlrc-bench: Can language agents solve machine learning research challenges? *arXiv preprint arXiv:2504.09702*, 2025.
- [912] Yuyu Zhang, Jing Su, Yifan Sun, Chenguang Xi, Xia Xiao, Shen Zheng, Anxiang Zhang, Kaibo Liu, Daoguang Zan, Tao Sun, et al. Seed-coder: Let the code model curate data for itself. *arXiv preprint arXiv:2506.03524*, 2025.
- [913] Zhihan Zhang, Yixin Cao, and Lizi Liao. Enhancing chart-to-code generation in multimodal large language models via iterative dual preference learning. *arXiv preprint arXiv:2504.02906*, 2025.
- [914] Chenlong Zhao, Xiwen Zhou, Xiaopeng Xie, and Yong Zhang. Hierarchical attention graph for scientific document summarization in global and local level. *arXiv preprint arXiv:2405.10202*, 2024.
- [915] Haiteng Zhao, Chang Ma, Fangzhi Xu, Lingpeng Kong, and Zhi-Hong Deng. Biomaze: Benchmarking and enhancing large language models for biological pathway reasoning. *arXiv preprint arXiv:2502.16660*, 2025.
- [916] Haiyan Zhao, Fan Yang, Bo Shen, Himabindu Lakkaraju, and Mengnan Du. Towards uncovering how large language model works: An explainability perspective. *arXiv preprint arXiv:2402.10688*, 2024.
- [917] Henry Hengyuan Zhao, Wenqi Pei, Yifei Tao, Haiyang Mei, and Mike Zheng Shou. Interfeedback: Unveiling interactive intelligence of large multimodal models via human feedback. *arXiv preprint arXiv:2502.15027*, 2025.
- [918] Penghai Zhao, Qinghua Xing, Kairan Dou, Jinyu Tian, Ying Tai, Jian Yang, Ming-Ming Cheng, and Xiang Li. From words to worth: Newborn article impact prediction with llm. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 1183–1191, Jan 2025. URL <https://ojs.aaai.org/index.php/AAAI/article/view/32106/34261>.
- [919] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 1(2), 2023.
- [920] Xuanle Zhao, Xianzhen Luo, Qi Shi, Chi Chen, Shuo Wang, Zhiyuan Liu, and Maosong Sun. Chart-coder: Advancing multimodal large language model for chart-to-code generation. *arXiv preprint arXiv:2501.06598*, 2025.

- [921] Xuanle Zhao, Zilin Sang, Yuxuan Li, Qi Shi, Shuo Wang, Duzhen Zhang, Xu Han, Zhiyuan Liu, and Maosong Sun. Autoreproduce: Automatic ai experiment reproduction with paper lineage. *arXiv preprint arXiv:2505.20662*, 2025.
- [922] Xueliang Zhao, Wenda Li, and Lingpeng Kong. Decomposing the enigma: Subgoal-based demonstration learning for formal theorem proving. *arXiv preprint arXiv:2305.16366*, 2023.
- [923] Yiming Zhao, Yongjia Zhao, Jian Wang, and Zhuo Wang. Artificial intelligence meets laboratory automation in discovery and synthesis of metal–organic frameworks: A review. *Industrial & Engineering Chemistry Research*, 64(9):4637–4668, Feb 2025.
- [924] Bo Zheng, Zhouyang Li, Fuxuan Wei, Qiguang Chen, Libo Qin, and Wanxiang Che. Hit-scir at mmnlu-22: Consistency regularization for multilingual spoken language understanding. *arXiv preprint arXiv:2301.02010*, 2023.
- [925] Kunhao Zheng, Jesse Michael Han, and Stanislas Polu. Minif2f: a cross-system benchmark for formal olympiad-level mathematics. *arXiv preprint arXiv:2109.00110*, 2021.
- [926] Mingyu Zheng, Xinwei Feng, Qingyi Si, Qiaoqiao She, Zheng Lin, Wenbin Jiang, and Weiping Wang. Multimodal table understanding. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9102–9124, Jun 2024.
- [927] Ruiqi Zheng, Liang Qu, Bin Cui, Yuhui Shi, and Hongzhi Yin. Automl for deep recommender systems: A survey. *ACM Trans. Inf. Syst.*, 41(4), March 2023. ISSN 1046-8188. doi: 10.1145/3579355. URL <https://doi.org/10.1145/3579355>.
- [928] Shenghe Zheng, Qianjia Cheng, Junchi Yao, Mengsong Wu, Ning Ding, Yu Cheng, Shuyue Hu, Lei Bai, Dongzhan Zhou, Ganqu Cui, et al. Scaling physical reasoning with the physics dataset. *arXiv preprint arXiv:2506.00022*, 2025.
- [929] Tianshi Zheng, Zheye Deng, Hong Ting Tsang, Weiqi Wang, Jiaxin Bai, Zihao Wang, and Yangqiu Song. From automation to autonomy: A survey on large language models in scientific discovery. *arXiv preprint arXiv:2505.13259*, 2025.
- [930] Xiaofeng Zheng and Jian Zhang. The usage of a transformer based and artificial intelligence driven multidimensional feedback system in english writing instruction. *Scientific Reports*, 15(1):1–22, Jun 2025.
- [931] Tianyang Zhong, Zhengliang Liu, Yi Pan, Yutong Zhang, Yifan Zhou, Shizhe Liang, Zihao Wu, Yanjun Lyu, Peng Shu, Xiaowei Yu, et al. Evaluation of openai o1: Opportunities and challenges of agi. *arXiv preprint arXiv:2409.18486*, 2024.
- [932] Junting Zhou, Wang Li, Yiyan Liao, Nengyuan Zhang, Tingjia Miaoand Zhihui Qi, Yuhan Wu, and Tong Yang. Academicbrowse: Benchmarking academic browse ability of llms. *arXiv preprint arXiv:2506.13784*, 2025.
- [933] Li Zhou, Ruijie Zhang, Xunlian Dai, Daniel Hershcovich, and Haizhou Li. Large language models penetration in scholarly writing and peer review. *arXiv preprint arXiv:2502.11193*, 2025.
- [934] Ruiyang Zhou, Lu Chen, and Kai Yu. Is llm a reliable reviewer? a comprehensive evaluation of llm on automatic paper reviewing tasks. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 9340–9351, May 2024.

- [935] Yangqiaoyu Zhou, Haokun Liu, Tejes Srivastava, Hongyuan Mei, and Chenhao Tan. Hypothesis generation with large language models. *arXiv preprint arXiv:2404.04326*, 2024.
- [936] Yilun Zhou, Caiming Xiong, Silvio Savarese, and Chien-Sheng Wu. Shared imagination: Llms hallucinate alike. *arXiv preprint arXiv:2407.16604*, 2024.
- [937] Haonan Zhu, Mary Silva, Jose Cadena, Braden Soper, Michał Lisicki, Braian Peetoom, Sergio E Baranzini, Shivshankar Sundaram, Priyadip Ray, and Jeff Drocco. Deep active learning based experimental design to uncover synergistic genetic interactions for host targeted therapeutics. *arXiv preprint arXiv:2502.01012*, 2025.
- [938] Kun Zhu, Xiaocheng Feng, Xiachong Feng, Yingsheng Wu, and Bing Qin. Hierarchical catalogue generation for literature review: a benchmark. *arXiv preprint arXiv:2304.03512*, 2023.
- [939] Kunlun Zhu, Jiaxun Zhang, Ziheng Qi, Nuoxing Shang, Zijia Liu, Peixuan Han, Yue Su, Haoifei Yu, and Jiaxuan You. Safescientist: Toward risk-aware scientific discoveries by llm agents. *arXiv preprint arXiv:2505.23559*, 2025.
- [940] Minjun Zhu, Yixuan Weng, Linyi Yang, and Yue Zhang. Deepreview: Improving llm-based paper review with human-like deep thinking process. *arXiv preprint arXiv:2503.08569*, 2025.
- [941] Minjun Zhu, Qiujie Xie, Yixuan Weng, Jian Wu, Zhen Lin, Linyi Yang, and Yue Zhang. Ai scientists fail without strong implementation capability. *arXiv preprint arXiv:2506.01372*, 2025.
- [942] Qing Zhu, Yan Huang, Donglai Zhou, Luyuan Zhao, Lulu Guo, Ruyu Yang, Zixu Sun, Man Luo, Fei Zhang, Hengyu Xiao, et al. Automated synthesis of oxygen-producing catalysts from martian meteorites by a robotic ai chemist. *Nature Synthesis*, 3(3):319–328, 2024.
- [943] Ronghang Zhu, Dongliang Guo, Daiqing Qi, Zhixuan Chu, Xiang Yu, and Sheng Li. A survey of trustworthy representation learning across domains. *ACM Transactions on Knowledge Discovery from Data*, 18(7):1–53, Jun 2024.
- [944] Xiaodan Zhu, Peter Turney, Daniel Lemire, and André Vellino. Measuring academic influence: Not all citations are equal. *Journal of the Association for Information Science and Technology*, 66(2):408–427, Jan 2015.
- [945] Jun Zhuang and Casey Kennington. Understanding survey paper taxonomy about large language models via graph representation learning. In Tirthankar Ghosal, Amanpreet Singh, Anita Waard, Philipp Mayr, Aakanksha Naik, Orion Weller, Yoonjoo Lee, Shannon Shen, and Yanxia Qin, editors, *Proceedings of the Fourth Workshop on Scholarly Document Processing (SDP 2024)*, pages 58–69, Bangkok, Thailand, August 2024. Association for Computational Linguistics. URL <https://aclanthology.org/2024.sdp-1.6/>.
- [946] Zhenzhen Zhuang, Jiandong Chen, Hongfeng Xu, Yuwen Jiang, and Jialiang Lin. Large language models for automated scholarly paper review: A survey. *arXiv preprint arXiv:2501.10326*, 2025.
- [947] Ziyu Zhuang, Qiguang Chen, Longxuan Ma, Mingda Li, Yi Han, Yushan Qian, Haopeng Bai, Zixian Feng, Weinan Zhang, and Ting Liu. Through the lens of core competency: Survey on evaluation of large language models. *arXiv preprint arXiv:2308.07902*, 2023.

- [948] Yoel Zimmermann, Adib Bazgir, Alexander Al-Feghali, Mehrad Ansari, L Catherine Brinson, Yuan Chiang, Defne Circi, Min-Hsueh Chiu, Nathan Daelman, Matthew L Evans, et al. 34 examples of llm applications in materials science and chemistry: Towards automation, assistants, agents, and accelerated scientific discovery. *arXiv preprint arXiv:2505.03049*, 2025.
- [949] James Zou and Eric J Topol. The rise of agentic ai teammates in medicine. *The Lancet*, 405(10477): 457, Feb 2025.
- [950] Apostolos Zournas, Matthew R Incha, Tijana Radivojevic, Vincent Blay, Jose Manuel Martí, Zak Costello, Matthias Schimdt, Tan Chung, Mitchell G Thompson, Allison Pearson, et al. Machine learning-led semi-automated medium optimization reveals salt as key for flaviolin production in *pseudomonas putida*. *Communications Biology*, 8(1):630, Apr 2025.