

Transforming Science with Large Language Models: A Survey on AI-assisted Scientific Discovery, Experimentation, Content Generation, and Evaluation

STEFFEN EGER, University of Technology Nuremberg (UTN), Germany

YONG CAO, University of Tübingen, Tübingen AI Center, Germany

JENNIFER D’SOUZA, TIB Leibniz Information Centre for Science and Technology, Germany

ANDREAS GEIGER, University of Tübingen, Tübingen AI Center, Germany

CHRISTIAN GREISINGER, University of Technology Nuremberg (UTN), Germany

STEPHANIE GROSS, Austrian Research Institute for Artificial Intelligence, Austria

YUFANG HOU, IT:U Interdisciplinary Transformation University Austria, Austria

BRIGITTE KRENN, Austrian Research Institute for Artificial Intelligence, Austria

ANNE LAUSCHER, University of Hamburg, Germany

YIZHI LI, University of Manchester, United Kingdom

CHENGHUA LIN, University of Manchester, United Kingdom

NAFISE SADAT MOOSAVI, University of Sheffield, United Kingdom

WEI ZHAO, University of Aberdeen, United Kingdom

TRISTAN MILLER, University of Manitoba, Canada

With the advent of large multimodal language models, science is now at a threshold of an AI-based technological transformation. Recently, a plethora of new AI models and tools have been proposed, promising to empower researchers and academics worldwide to conduct their research more effectively and efficiently. This includes all aspects of the research cycle, especially (1) searching for relevant literature; (2) generating research ideas and conducting experimentation; generating (3) text-based and (4) multimodal content (e.g., scientific figures and diagrams); and (5) AI-based automatic peer review. In this survey, we provide an in-depth overview over these recent advances, which promise to fundamentally alter the scientific research process for good. Our survey covers the five aspects outlined above, indicating relevant datasets, methods and results (including evaluation) as well as limitations and scope for future research. Ethical concerns regarding shortcomings of these tools and potential for misuse (fake science, plagiarism, harms to research integrity) take a particularly prominent place in our discussion. We hope that our survey will not only become a reference guide for newcomers to the field but also a catalyst for new AI-based initiatives in the area of “AI4Science”.

CCS Concepts: • **Social and professional topics** → *Assistive technologies*; • **Applied computing** → *Physical sciences and engineering; Life and medical sciences; Law, social and behavioral sciences*; • **Computing methodologies** → **Natural language processing; Artificial intelligence**; • **General and reference** → **Surveys and overviews**.

Authors’ Contact Information: [Steffen Eger](mailto:steffen.eger@utn.de), steffen.eger@utn.de, University of Technology Nuremberg (UTN), Nuremberg, Germany; [Yong Cao](mailto:yong.cao@uni-tuebingen.de), yong.cao@uni-tuebingen.de, University of Tübingen, Tübingen AI Center, Tübingen, Germany; [Jennifer D’Souza](mailto:jennifer.dsouza@tib.eu), jennifer.dsouza@tib.eu, TIB Leibniz Information Centre for Science and Technology, Hannover, Germany; [Andreas Geiger](mailto:a.geiger@uni-tuebingen.de), a.geiger@uni-tuebingen.de, University of Tübingen, Tübingen AI Center, Tübingen, Germany; [Christian Greisinger](mailto:christian.greisinger@utn.de), christian.greisinger@utn.de, University of Technology Nuremberg (UTN), Nuremberg, Germany; [Stephanie Gross](mailto:stephanie.gross@ofai.at), stephanie.gross@ofai.at, Austrian Research Institute for Artificial Intelligence, Vienna, Austria; [Yufang Hou](mailto:yufang.hou@it-u.at), yufang.hou@it-u.at, IT:U Interdisciplinary Transformation University Austria, Linz, Austria; [Brigitte Krenn](mailto:brigitte.krenn@ofai.at), brigitte.krenn@ofai.at, Austrian Research Institute for Artificial Intelligence, Vienna, Austria; [Anne Lauscher](mailto:anne.lauscher@uni-hamburg.de), anne.lauscher@uni-hamburg.de, University of Hamburg, Hamburg, Germany; [Yizhi Li](mailto:yizhi.li-2@manchester.ac.uk), yizhi.li-2@manchester.ac.uk, University of Manchester, Manchester, United Kingdom; [Chenghua Lin](mailto:chenghua.lin@manchester.ac.uk), chenghua.lin@manchester.ac.uk, University of Manchester, Manchester, United Kingdom; [Nafise Sadat Moosavi](mailto:n.s.moosavi@sheffield.ac.uk), n.s.moosavi@sheffield.ac.uk, University of Sheffield, Sheffield, United Kingdom; [Wei Zhao](mailto:wei.zhao@abdn.ac.uk), wei.zhao@abdn.ac.uk, University of Aberdeen, Aberdeen, United Kingdom; [Tristan Miller](mailto:Tristan.Miller@umanitoba.ca), Tristan.Miller@umanitoba.ca, University of Manitoba, Winnipeg, Manitoba, Canada.

Additional Key Words and Phrases: Language Language Models, Science, AI4Science, Search, Experimentation, Idea Generation, Multimodal Content Generation, Evaluation, Peer Review

1 Introduction

Throughout history, science has undergone a number of paradigm shifts, culminating in today's era of data-intensive exploration [88]. Although new tools and frameworks have accelerated the pace of scientific discovery, its basic steps have remained unchanged for centuries. These include (1) conception of a research question or problem, typically arising from a gap in disseminated knowledge; (2) collection and study of existing literature or data relevant to the problem; (3) formulation of a falsifiable hypothesis; (4) design and execution of experiments to test this hypothesis; (5) analysis and interpretation of the resulting data; and (6) reporting on the findings, allowing for their exploitation in real-world applications or as a source of knowledge for a further iteration of the scientific cycle. A more detailed discussion of these steps is provided in Appendix A.1.

With the advent of large multimodal foundation models such as ChatGPT, Gemini, Qwen, or DeepSeek, many research fields and sectors of everyday life now stand at the threshold of an AI-based technological transformation. Science is no exception. A recent study analyzed approximately 148,000 papers from 22 non-CS fields that cited large language models (LLMs) between 2018 and 2024, reporting a growing prevalence of LLMs in these disciplines [176]. Additionally, a very recent survey among almost 5000 researchers in more than 70 countries, by the American Publishing Company Wiley, suggests that AI adoption is embraced by a majority of researchers who think that AI will become mainstream in science within the next two years, despite current usages often limited to forms of writing assistance.¹

While science has traditionally relied on human ingenuity and labor for generating research ideas, formulating hypotheses, searching for relevant literature, conducting experiments, and reporting results, recent advancements in AI have introduced a surge of models and tools promising to assist researchers at every stage of this cycle. These include models like Elicit or ORKG ASK for search; models like The AI Scientist [151] for experimentation; and models like AutomaTikZ [14] and DeTikZify [15] for multimodal scientific content generation. Moreover, there is even research investigating the extent to which these AI models can evaluate scientific outcomes through automated peer review [259]. These AI-driven advancements promise to accelerate the scientific process, leading to unexpected discoveries, improved documentation, and more accessible research communication.²

Despite this rapid progress, to our knowledge, there is no comprehensive survey covering the full breadth of AI-assisted tools, models, and functionalities available for supporting and improving the research cycle. Existing reviews are typically domain-specific, such as in the social sciences [249] or branches of physics [265].³ To address this gap, our survey provides an extensive overview of five central aspects of the research cycle where AI is making transformative contributions: (1) search and content summarization (Section 3.1); (2) scientific experimentation and research idea generation (Section 3.2); (3) unimodal content generation, such as drafting titles, abstracts, suggesting citations, and assisting in text refinement (Section 3.3); (4) multimodal content generation, including the creation and interpretation of figures, tables, slides, and posters (Section 3.4); and (5) AI-assisted peer review processes (Section 3.5). The recent Wiley survey underscores the significance of this endeavor, highlighting that “63% [of respondents indicated] a lack of clear guidelines and consensus on what uses of AI are accepted in their field and/or the need for more training and skills.”

¹<https://www.wiley.com/en-us/ai-study>, <https://www.nature.com/articles/d41586-025-00343-5>

²The benefits are expected to be particularly significant for non-native English speakers and those with lower technical skills, potentially increasing diversity and inclusivity in research.

³We note two contemporaneous works developed completely independently from us [153, 266]. Both are substantially narrower in scope than this survey; for example, Luo et al. [153] neither cover multimodal approaches to scientific content synthesis nor search and also do not address ethical concerns in nearly the same depth as we do.

When it comes to the use of AI tools in scientific research, ethical considerations are paramount. These tools exhibit various limitations, including (i) hallucinating and fabricating content, (ii) exhibiting bias, (iii) having limited reasoning abilities, (iv) lacking proper evaluation mechanisms, and (v) posing significant environmental costs. Additional concerns include risks of fake science, plagiarism, and diminished human oversight in scientific processes. The European Union has recently released guidelines on the responsible use of AI in science, emphasizing that, while [r]esearch is one of the sectors that could be most significantly disrupted by generative AI,” and AI has great potential for accelerating scientific discovery and improving the effectiveness and pace of research and verification processes,” it also “raises questions about the ability of current models to combat deceptive scientific practices and misinformation.”⁴ In this survey, we address these ethical considerations by including dedicated discussions within each of the five aspects of the research cycle and a standalone discussion in Section 4.

As shown in Fig. 1, the rest of this paper is structured as follows: §2 discusses the methodological approach of our survey. In §3, each subsection describes the application of AI to an individual scientific task (literature search, experimental design, writing, etc.). For each task, we discuss important datasets, state-of-the-art methods and results, ethical concerns, domains of application, limitations, and future directions. §4 presents ethical considerations beyond individual tasks. Finally, §5 summarizes the benefits and challenges of AI in science.

Resources related to this survey are available at <https://github.com/NL2G/TransformingScienceLLMs>. A list of abbreviations and further material is available in our appendix.

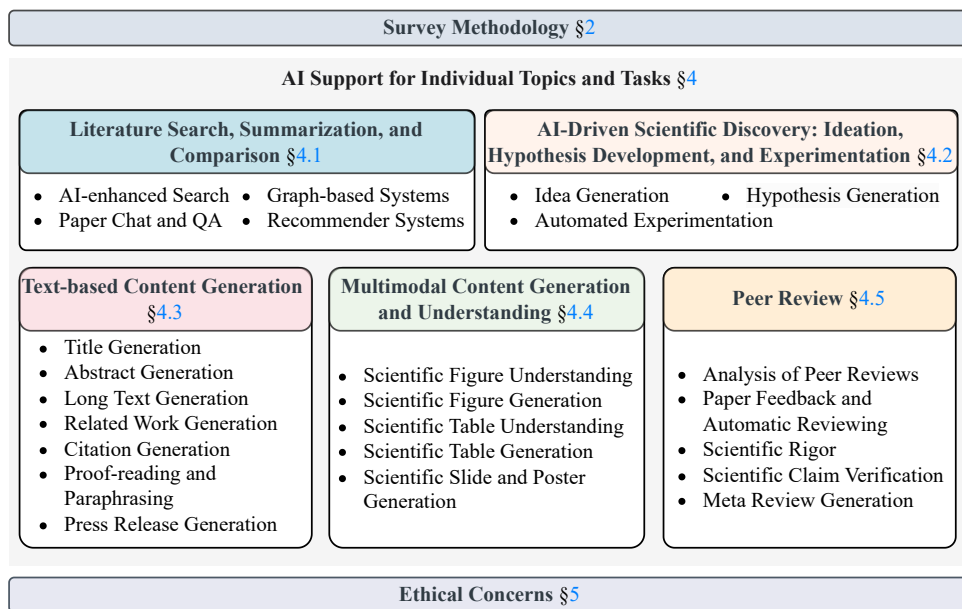


Fig. 1. Overview of the survey structure, including our survey methodology, five AI-assisted topics or tasks, and ethical considerations.

⁴https://research-and-innovation.ec.europa.eu/document/download/2b6cf7e5-36ac-41cb-aab5-0d32050143dc_en?filename=ec_rtd_ai-guidelines.pdf

2 Survey Methodology

This article offers a detailed, disciplinarily contextualized survey of state-of-the-art AI applications in scientific research, spanning every stage from the initial conception of ideas to the dissemination of results. It is intended primarily to help researchers in fields within AI (natural language processing (NLP), computer vision (CV), etc.) quickly familiarize themselves with the transdisciplinary foundations of and latest developments in this broad-ranging and rapidly evolving research area. Some of the material will also be useful to policymakers, practitioners, and research collaborators in adjacent fields, including human–computer interaction, library and information science, communication studies, metascience, science journalism, and research ethics.

We believe our contribution to be timely because, despite the growing interest in the topic, its researchers are only just now coalescing into a community with dedicated dissemination venues. Recent examples include the workshops Natural Scientific Language Processing and Research Knowledge Graphs (NSLP) [189], Foundation Models for Science (FM4Science), AI & Scientific Discovery (AISD), and Towards a Knowledge-grounded Scientific Research Lifecycle (AI4Research), all of which held their first editions in 2024 or 2025. The few existing reviews of AI-for-science literature have addressed only isolated topics or application domains. The earliest examples (e.g., [58, 122]), now long out of date, tend to be organized into case studies of AI for specialized tasks such as equation or drug discovery. More recent surveys, such as [85], cover a wider variety of application domains but focus on a narrower sector of the scientific life cycle, and are pitched more towards the potential users of the AI tools than towards AI researchers aiming to understand and advance the underlying data sets, methodologies, and evaluation metrics.

Given our topic’s wide scope, rapid progress, and dependence on knowledge and methods from different domains, we have opted to take a *narrative approach* to our survey. This methodology allows for greater freedom in the selection and framing of relevant papers [111], which promotes “breadth of literature coverage and flexibility to deal with evolving knowledge and concepts” [21] as well as the ability to “bridge related areas of work, provoke thoughts, inspire new theoretical models, and direct future efforts in a research domain” [171]. *Systematic reviews*, while regarded as more objective, are better suited to relatively narrow topics with well-defined, empirical research questions [171]. Accordingly, we have adopted no fixed inclusion or exclusion criteria for the studies referenced in this survey, but have rather selected them on the basis of our own relevance judgements. In assembling the co-authors for this survey, we have therefore endeavoured to include researchers actively publishing in each of the various subtopics we cover.

3 AI Support for Individual Topics and Tasks

3.1 Literature Search, Summarization, and Comparison

The rapid growth of scientific literature presents a significant challenge for researchers who need to search, analyze, and summarize vast amounts of information efficiently. AI-powered tools are transforming these tasks by leveraging NLP, machine learning (ML), LLMs, citation and knowledge graphs (KGs) to automate the retrieval, extraction, and summarization of scientific information. Unlike traditional search engines that rely on basic keyword matching, AI-equipped systems provide context-aware, semantic search with additional features that enhance the overall search experience. These systems go beyond finding relevant papers; they generate answers to research questions from the search results, provide structured summaries, and offer comparative insights, helping researchers identify gaps, trends, and contradictions across multiple studies.

3.1.1 Data. Scientific search engines rely on vast publisher databases to provide access to scientific literature. Understanding the classification of these repositories is essential for assessing search engines’ coverage, reliability, and effectiveness in evidence-based research. Repositories vary by **access model**, **subject focus**, and **content type**, each serving a distinct role in academic discovery and knowledge dissemination. By **access model**, repositories fall into *open access repositories*, which provide unrestricted access to research articles (e.g., [PubMed Central](#), [arXiv](#)); *subscription-based repositories*, requiring institutional or individual subscriptions (e.g., [ScienceDirect](#), [SpringerLink](#)); and *hybrid repositories*, offering both free and paywalled content (e.g., [Taylor & Francis Online](#), [Oxford Academic](#)). By **subject focus**, repositories are either *multidisciplinary*, covering broad disciplines (e.g., [Web of Science](#), [Scopus](#)), or *subject-specific*, specializing in fields such as medicine ([PubMed](#)), physics ([INSPIRE-HEP](#)), and social sciences ([SSRN](#)). By **content type**, *institutional repositories* archive research outputs from specific organizations (e.g., [MIT DSpace](#), [Harvard DASH](#)); *preprint repositories* enable early dissemination of research before peer review (e.g., [bioRxiv](#), [chemRxiv](#)); and *government and public sector repositories* provide access to publicly funded research (e.g., [NASA ADS](#), [OpenAIRE](#)). *Data repositories* (e.g., [Dryad](#), [Zenodo](#)) store research datasets, supporting transparency and reproducibility, while *aggregator repositories* (e.g., [BASE](#), [CORE](#) [113]) index content from multiple sources for broader searches. Lastly, *grey literature repositories* (e.g., [OpenGrey](#), [EThOS](#)) provide access to non-traditional research outputs such as theses, reports, and white papers, which may not be available through conventional publisher platforms.

The structure of scientific repositories shapes AI-enhanced search. While broad AI-based search engines like [Elicit](#) and [ORKG ASK](#) query multiple publisher repositories, similar to [Google Scholar](#), tools like [NotebookLM](#) focus on user-selected documents, and recommender systems such as [Scholar Inbox](#) rank new literature by relevance. AI-driven search enables customizable knowledge bases while optimizing discovery, retrieval, and personalization in research.

3.1.2 Methods and Results. This section surveys state-of-the-art AI-enhanced scientific discovery tools, identified as four main types based on their core functionality: (1) **AI-enhanced search**, which retrieves relevant literature from vast repositories; (2) **graph-based systems**, which map relationships between research concepts and publications; (3) **paper chat and QA**, which enable interactive exploration of scientific content; and (4) **recommender systems**, which suggest relevant papers based on user preferences. A detailed overview of these tools is provided in Table 1. Additionally, two traditional scientific discovery tools, namely search engines and benchmarks with leaderboards, are discussed in Appendix A.2.1.

AI-enhanced Search. Platforms such as [Elicit](#), [Consensus](#), [OpenScholar](#), and [SciSpace](#) leverage AI, including LLMs, to extend beyond traditional search by enabling semantic search, paper summarization, evidence synthesis, and trend analysis. Unlike conventional search engines that rely on keyword matching, these tools use NLP and machine learning to extract key insights, synthesize information to answer research queries [76], and generate structured summaries. Their ability to quickly summarize and categorize findings—such as study outcomes, methodologies, and limitations—helps researchers efficiently compare and interpret literature.

Graph-based Systems. Graph-based systems such as [ORKG ASK](#) are designed to facilitate structured access to scientific knowledge. Unlike conventional paper search engines, they leverage a KG that organizes research contributions as structured data rather than unstructured text. Such contributions are typically extracted from the abstract, introduction, and result sections [56, 175]. Those systems enable users to ask complex, domain-specific questions and receive answers synthesized from semantically structured scientific data. They typically use techniques such as KG-based reasoning and retrieval-augmented generation (RAG) to extract relevant information from the KG, providing more interpretable and

Platform	Search	Recommendations	Citation Analysis	Trending Analysis	Author Profiles	Visualization Tools	Paper Chat	Idea Generation	Paper Writing	Summarization	Paper Review	Datasets	Code Repositories	LLM Integration	Web API	Personalization	Cost	Data Source
AI-Enhanced Search	Elicit	✓				✓	✓		✓	✓		✓					Freemium	125 million
	OpenScholar	✓	✓			✓			✓			✓					Free	45 million
	Undermind	✓	✓			✓				✓		✓		✓			Premium	over 200 million
	Perplexity	✓				✓	✓		✓	✓		✓					Freemium	
	Consensus	✓	✓			✓			✓			✓	✓				Freemium	over 200 million
	SciSpace	✓	✓			✓	✓		✓	✓			✓				Freemium	
	scienceQA	✓	✓	✓		✓	✓		✓	✓			✓				Freemium	220 million
	PaperQA2					✓					✓	✓					Free	
	Paperguide	✓	✓			✓			✓	✓			✓				Freemium	
	HyperWrite	✓				✓	✓	✓	✓	✓			✓				Premium	
	ResearchKick	✓				✓	✓	✓	✓	✓			✓		✓		Premium	
Graph-Based	Connected Papers	✓	✓		✓												Freemium	214 million
	ScholarGPS	✓		✓	✓	✓	✓										Free	over 200 million
	CiteSpace				✓		✓										Freemium	
	Sci2						✓										Free	
	NLP KG	✓	✓	✓		✓											Free	
	ORKG ASK	✓	✓						✓				✓				Free	76 million
Paper Chat	ChatGPT	✓				✓	✓	✓	✓	✓		✓	✓				Freemium	10 pdf files
	Claude	✓				✓	✓	✓	✓	✓		✓	✓				Freemium	5 pdf files
	Deepseek	✓				✓	✓	✓	✓	✓		✓	✓				Free	
	Research		✓			✓	✓		✓	✓		✓					Freemium	1 pdf file
	NotebookLM					✓	✓		✓	✓		✓			✓		Freemium	50 pdf files
	Enago Read	✓	✓			✓	✓		✓	✓		✓			✓		Freemium	1 pdf file
	DocAnalyzer.AI			✓		✓	✓		✓	✓		✓	✓	✓			Premium	few pdf files
	CoralAI		✓			✓	✓		✓	✓		✓					Freemium	1 pdf file
	ExplainPaper					✓	✓		✓	✓		✓					Freemium	1 pdf file
	ChatPDF	✓	✓			✓	✓		✓	✓		✓					Premium	1 pdf file
Recommender	Arxiv Sanity	✓	✓	✓											✓		Free	
	Scholar Inbox	✓	✓	✓									✓		✓		Free	
	ResearchTrend.ai	✓			✓												Freemium	
	TrendingPapers	✓	✓		✓												Free	
	Bytez	✓			✓				✓	✓			✓	✓			Freemium	
	Notesum.ai	✓	✓	✓					✓				✓		✓		Freemium	
	Research Rabbit	✓	✓			✓											Free	

Table 1. Overview of popular literature search, summarization and comparison tools and their key features.

verifiable answers compared to traditional LLM-based QA systems. [CiteSpace](#) and [Sci2](#) are specialized bibliometric analysis and network analysis tools to study the structure and evolution of scientific research. [CiteSpace](#) focuses on identifying research trends, keyword co-occurrence networks, and citation bursts, using visual analytics to highlight emerging topics and influential papers using graphs. [Sci2](#) is a more general-purpose tool designed for analyzing scholarly datasets, enabling users to perform network analysis, geospatial mapping, and temporal modeling of scientific literature and collaboration patterns. [Connected Papers](#) is a scientific literature exploration tool designed to help researchers discover related papers based on a given seed paper. Unlike traditional citation-based systems, it builds a graph of papers using a similarity metric derived from co-citation and bibliographic coupling analysis. The platform constructs a network where each node represents a paper, and edges indicate similarity based on shared references and citations rather than direct citation links. This approach allows users to find relevant papers that may not be directly cited but

are conceptually related. Graph-based visualizations provide an intuitive way to explore clusters of research, identify foundational and emerging works, and track the evolution of scientific ideas.

Paper Chat and QA. Paper chat and question-answering (QA) systems such as [ChatGPT](#), [Deepseek Chat](#), [NotebookLM](#), [ExplainPaper](#), [ChatPDF](#), and [DocAnalyzer.AI](#) allow users to interact with scientific papers by asking questions and receiving responses based on the document's content. They typically process a limited number of user-provided PDFs or text from specific websites. The core technology behind them is RAG [8, 105, 126], a technique that combines information retrieval with LLMs to improve accuracy and grounding. A typical RAG system first partitions the document into smaller sections and converts them into vector representations using embedding models. Upon a user query, the system retrieves the most relevant sections based on semantic similarity and passes them as context to an LLM, which then generates a response. This mechanism ensures that answers are directly grounded in the provided documents rather than relying solely on the model's pre-trained knowledge, enhancing factual reliability and interpretability. Some systems incorporate LLM agents [22, 138, 223] that can reason over retrieved information, summarize findings, or extract key insights. These agents can follow multi-step reasoning strategies to provide more nuanced responses, such as synthesizing information from multiple sections or explaining technical terms in simpler language. By anchoring responses to document content, RAG-based systems mitigate hallucinations and make it easier for users to verify claims by checking the referenced passages. The effectiveness of these systems depends on the quality of document chunking, the efficiency of retrieval, and the model's ability to integrate information into coherent, context-aware answers.

Recommender Systems. Scientific paper recommender systems such as [Arxiv Sanity](#), [Scholar Inbox](#), [ResearchTrend.ai](#), and [Research Rabbit](#) leverage machine learning and information retrieval techniques to help researchers discover relevant literature. These systems generally fall into two main categories: content-based filtering, collaborative filtering and hybrid approaches. Content-based methods [4, 16] analyze the text of papers to build representations that capture their meaning. Traditional approaches rely on sparse abstract or document representations such as TF-IDF [210], which assigns importance to words based on their frequency and distinctiveness in a corpus. More advanced models use dense abstract or document embeddings derived from neural networks, such as SPECTER [43] or GTE [139], which map papers into a high-dimensional vector space where similar documents are close to each other. The Massive Text Embedding Benchmark (MTEB) [164] ranks many state-of-the-art embedding models on a comprehensive benchmark comprising various different datasets and tasks. These embeddings enable fast similarity searches and improve over simple keyword matching. In contrast, collaborative filtering [12, 236] relies on user interactions, such as downloads, bookmarks, and citations, to recommend papers based on the behavior of similar users. One challenge of pure collaborative filtering is the cold start problem, where new papers or users lack sufficient data for recommendations. To mitigate this, many modern systems employ hybrid approaches, such as two-tower architectures [46, 255, 258]. These models learn separate representations for papers and users, combining textual embeddings with user interaction data to generate more personalized recommendations. State-of-the-art systems often use a mix of these techniques to balance relevance, novelty, and diversity. The effectiveness of these systems depends on the quality of embeddings, the availability of interaction data, and the efficiency of ranking algorithms that surface the most useful papers.

3.1.3 Ethical Concerns. The use of AI in scientific search, summarization and comparison raises ethical considerations, particularly in ensuring transparency, accountability, and equity. AI can significantly accelerate the pace of discovery, automate search tasks, and uncover patterns that may elude human researchers, but it also introduces risks and biases. Existing dynamics such as the Matthew effect, where well-known researchers receive disproportionate attention, might

be reinforced by the AI algorithms, intensifying inequalities. We believe that research should follow a human-centric approach, in which the human researcher is provided with advanced tools but remains fully responsible for executing the research and summarizing the results in research papers. It is also important to develop algorithms to reduce biases by recommending relevant work to researchers based on the *content* of the research, independent of the popularity of the authors. Tools that are able to uncover gaps in the existing literature might even lead to a more uniform allocation of researchers to topics, reducing the bias towards overpopulated areas.

3.1.4 Domains of Application. The search, summarization and comparison tools discussed in this section are general and apply to all fields of science. The presented benchmarks, however, are specific to the field of computer science in general and artificial intelligence in particular.

3.1.5 Limitations and Future Directions. Despite the significant advancements in AI-powered scholarly search systems, several limitations persist that hinder their full potential. One of the primary challenges is *data quality and coverage gaps*, as these systems often struggle with handling incomplete, non-standardized, or outdated data sources, which can lead to inaccuracies and inconsistencies in retrieved information. Additionally, *bias in AI models* remains a critical concern, where search and ranking algorithms may introduce biases based on training data, potentially influencing the visibility of certain research areas and limiting the diversity of perspectives presented to users. Another major limitation lies in *scalability and real-time processing*, as efficiently handling large-scale datasets while maintaining low latency and high retrieval accuracy remains a technical challenge. Addressing these limitations opens several promising future directions. One potential avenue is *enhanced personalization* which can be achieved by adapting search engines to user preferences, providing more tailored recommendations based on research interests and behavioral patterns. Lastly, fostering *interdisciplinary collaboration* through the integration of AI-powered search systems with other digital tools, such as data visualization platforms and research management software, could facilitate more comprehensive and insightful research outcomes. Addressing these challenges and exploring future directions will be crucial for realizing the full potential of AI-driven scholarly search and synthesis.

3.2 AI-Driven Scientific Discovery: Ideation, Hypothesis Generation, and Experimentation

Idea formation, hypothesis generation, and experimentation are fundamental to scientific discoveries. Idea formation, particularly in AI, focuses on proposing new tools or benchmarking existing ones. Hypothesis formation involves formulating specific, testable questions that guide empirical or theoretical justifications. Experimentation then tests hypotheses and evaluates ideas through systematic observation, data collection, and analysis. In AI research, this often includes benchmarking models, running simulations, or conducting ablation studies. Traditionally, these processes have been carried out by human researchers. However, in an age of rapidly growing scientific literature, efforts of moving from literature review to hypothesis and idea formation have become increasingly time-consuming. Experimentation adds further complexity, requiring careful methodological design, large-scale simulations, and in-depth result analysis. As AI-driven approaches accelerate hypothesis generation, idea formation, and evaluation, it is essential for the research community to assess a broad range of candidates, selecting those that are most meaningful, relevant, and potentially novel for further validation.

Dataset	Source	Data Size	Domain	Time Span	Task
SciMON [26]	ACL Anthology	135,814 papers	NLP	1952-2022	Idea Generation
IDEA Challenge [59]	University of Bristol	240 prototypes	Engineering	2022	Idea Generation
SPACE-IDEAS+ [71]	COLING	1020 ideas	Physics	2024	Idea Generation
TOMATO-Chem [253]	Nature and Science	51 papers	Chemistry	2024	Hypothesis Generation
TOMATO [252]	Unspecified	50 papers	Social Science	2023-2024	Hypothesis Generation
LLM4BioHypoGen [182]	PubMed	2,900 papers	Medicine	2000-2024	Hypothesis Generation
DevAI [271]	Meta	55 tasks	AI	2024	Automated Experimentation
ScienceAgentBench [35]	OSU NLP	44 papers	Diverse ⁵	2024	Automated Experimentation
SWE-bench [102]	ICLR	2,294 issues	SWE	2024	Automated Experimentation
MLGym-Bench [166]	Meta	13 tasks	Diverse ⁶	2025	Automated Experimentation

Table 2. Overview of datasets for idea and hypothesis generation and experimentation

3.2.1 Data. We survey diverse datasets for evaluating LLMs in hypothesis generation, idea formation and experimentation. These datasets are collected from various sources such as ACL Anthology, Nature and Science, PubMed and ICLR, and span multiple domains over NLP, Engineering, Physics, Chemistry and many more. For **idea generation**, datasets consist of paper abstracts [26], prototype ideas stemming from a design hackathon [59], innovative ideas for space-related technologies and missions [71]. For **hypothesis generation**, datasets contain scientific papers, and each paper is annotated by domain experts with background, research questions, related works and hypotheses in various domains [182, 252, 253]. For **automated experimentation**, datasets consist of structured problem-solving tasks that require generating or modifying Python code. These include AI development tasks with hierarchical user requirements [271], scientific discovery automation based on data and expert knowledge [35], software engineering tasks requiring code edits to resolve issues [102], and the MLGym framework [166], which evaluates agents on open-ended AI research tasks across multiple domains. Details of these datasets are presented in Table 2.

3.2.2 Methods and Results. Here, we discuss state-of-the-art methods and results in hypotheses generation, idea formation, and automated experimentation. Fig. 2 provides some examples for each approach.

Idea Generation. Research ideation has become a critical area where LLMs are increasingly applied to enhance novelty and accelerate discovery. Several methods have been proposed to improve the creative abilities of LLMs, focusing on iterative refinement, multi-agent systems, human alignment and evaluation [92, 131, 151, 184, 215].

For **iterative refinement**, Hu et al. [92] introduce an iterative planning and search framework aimed at enhancing the novelty and diversity of ideas generated by LLMs. By systematically retrieving external knowledge, the approach addresses the limitations of existing models in producing simplistic or repetitive suggestions. Similarly, Pu et al. [178] focus on iterative refinement by providing literature-grounded feedback. Representing research ideas as nodes on a canvas, their approach, IdeaSynth, facilitates the iterative exploration and composition of idea facets, enabling to develop more detailed and diverse ideas, particularly at various stages of ideation. For **human alignment**, other works seek to organize information in ways that mirror human research processes. Chain of Ideas (CoI) [131] proposes structuring literature into a chain to emulate the progressive development of research domains. This facilitates the identification of meaningful directions and has been shown to outperform existing methods in generating ideas comparable in quality to those produced by human researchers. Scideator [184], in contrast, focuses on recombining facets (e.g., purposes, mechanisms, and evaluations) from existing research papers to synthesize novel ideas. By incorporating automated novelty assessments, Scideator enables users to identify overlaps and refine their ideas. Moreover, research has explored fully **autonomous and multi-agent systems** for scientific discovery. For instance, the AI Scientist [151] presents a framework for automating the entire research pipeline, including idea generation, experiment execution, and paper

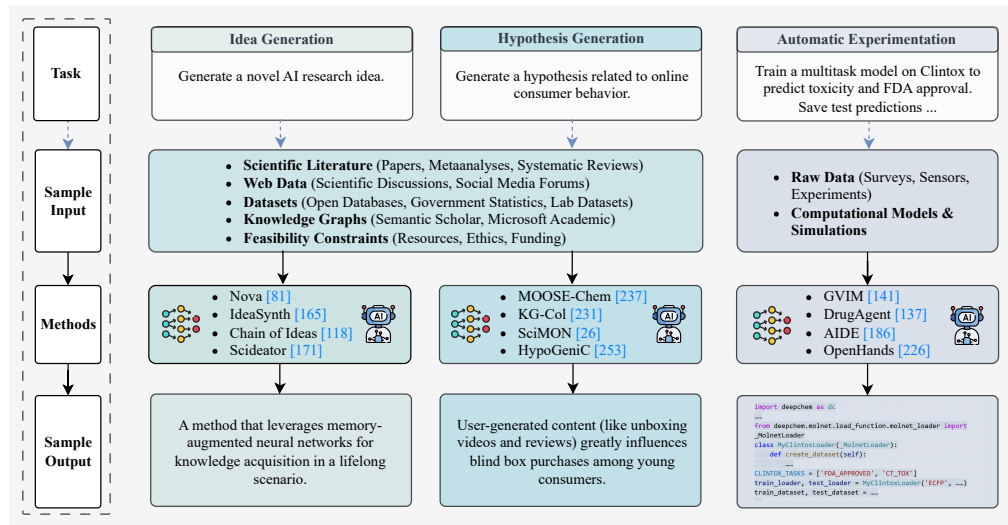


Fig. 2. Examples of idea generation, hypothesis generation, and automated experimentation follow a four-component structure: *task*, *sample input*, *methods*, and *sample output*. The *task* defines the goal of each process. *Sample input* consist of benchmark datasets for each task. *Methods* encompass relevant scientific approaches. *Sample output* differs by process: idea generation and hypothesis yield textual outputs (descriptions or explanations), whereas automated experimentation produces executable code.

writing. VirSci [215] employs a multi-agent system of virtual agents to collectively generate, evaluate, and refine ideas, which outperforms individual LLMs, underscoring the potential of teamwork in enhancing scientific innovation.

Hypothesis Generation. Recently, there have been many works that leverage LLMs to generate hypotheses and ideas [26, 147, 182, 235, 247, 253]. These works differ in their use of LLMs for addressing various technique challenges, including (a) handling long context input due to the need for LLMs to analyze related works, (b) strategies for refining LLMs to generate meaningful hypotheses, (c) lowering the possibility of generating hallucinated hypotheses and ideas.

For **hallucinated hypotheses**, Yang et al. [253] address the hallucination issue through a pipeline that starts with a search for related works. The identified related works and a given research question are provided as input for LLMs to generate hypotheses. The generated hypotheses are evaluated against ground-truth hypotheses published in Nature and Science. Their results show that many generated hypotheses exhibit a very high degree of similarity to the ground-truth ones. Regarding **long context input**, Chai et al. [26] focus on the efficient use of limited context size of LLMs. They introduce a selection mechanism that extracts important and relevant information from the literature and takes them as input for LLMs to generate hypotheses. Their results show that filtering out unnecessary information helps improve the quality of generated hypotheses. For **refinement strategies**, many works have explored strategies for refining LLMs to generate hypotheses. Major strategies include (a) few-shot learning, (b) fine-tuning on training data and (c) iterative refinement. For instance, Qi et al. [181] show that hypotheses generated by few-shot learning are judged by humans more testable than those generated in the zero-shot setup; while fine-tuning improves the overall quality of hypotheses, the improvement is limited to the domain of training data; in unseen domains, fine-tuning harms hypothesis quality, particularly the novelty aspect. Zhou et al. [270] iteratively refine hypotheses through reinforcement learning, with the aim of increasing the similarity between a given research problem and a generated hypothesis. A recent advancement in

this space is the AI co-scientist⁷, a multi-agent system that employs a generate-debate-evolve framework. It iteratively enhances its hypotheses through collaboration, grounding in prior evidence and tournament-based selection.

Automated Experimentation. Experimentation is central to AI-driven research, encompassing task formulation, implementation, evaluation, and iteration. Automated experimentation aims to streamline this workflow, with approaches like Neural Architecture Search [62] and AutoML [86]. LLMs further enhance this by enabling automation through natural language prompts. AutoML-GPT [227] and MLcopilot [264] use LLMs for hyperparameter tuning, while MLAGentBench [98] benchmarks fundamental automation tasks. Recent work explores advanced frameworks incorporating multi-agent collaboration, tree search, and iterative refinement for scientific experimentation.

For **multi-agent workflow**, GVIM [154] enhances chemical research with domain-specific functions, while DrugAgent [150] employs LLMs for task planning in drug discovery. AutoML-Agent [226] integrates retrieval-augmented planning for AutoML tasks, and MLAGentBench [98] benchmarks LLM-driven agents in machine learning experimentation. The Agent-as-a-Judge framework [271] introduces structured agent evaluation. For **tree search**, AIDE [201] applies Solution Space Tree Search to refine solutions in Kaggle challenges. The "Tree Search for Language Model Agents" framework [115] enables LLM agents to plan multi-step interactions using best-first tree search, pruning less promising options. SELA [39] combines LLM-generated insights with Monte Carlo Tree Search, iteratively refining machine learning experiments by selecting promising configurations and executing them. For **Iterative refinement**, APEX [44] automates LLM-based experimentation with an orchestrator, execution engine, benchmark generator, and model library. OpenHands [241] enables AI agents to interact with software, execute actions in a sandboxed runtime, and collaborate across tasks using predefined benchmarks.

Evaluation. While LLMs accelerate ideation and hypothesis generation, it is crucial for the research community to assess their usefulness and identify those worth further validation. Many evaluation approaches have been proposed, differing primarily in whether gold hypotheses and ideas are available. When **gold hypotheses and ideas are available**, metrics such as BLEU [168] have been used to assess hypothesis and idea quality by measuring similarity to known scientific discoveries [26, 253]. Qi et al. [182] leverage LLMs-as-a-metric to evaluate hypotheses based on four scientific aspects: (a) novelty, (b) relevance to the given background, (c) significance within the research community, and (d) verifiability, i.e., testability. When **gold hypotheses and ideas are unavailable**, a generated hypothesis or idea is typically assessed by human experts, who compare it with the given research question and provide feedback [252].

3.2.3 Ethical Concerns. In the area of idea generation, there is a risk of reinforcing established research paradigms. AI systems trained on the basis of existing literature may favor popular paths and neglect underrepresented research directions. As a result, unconventional ideas may be unintentionally marginalized. AI-generated hypotheses may lack transparency, making it difficult to assess their validity or underlying assumptions, which could lead to flawed experiments. For example, an AI might identify a statistical correlation in its training data and propose hypotheses without clearly revealing the underlying assumptions or data sources, making it difficult for researchers to verify its scientific soundness or hold anyone accountable if the hypotheses proves misleading. Automated experimentation presents its own ethical challenges. The speed and volume in which AI can design and execute experiments can lead to insufficient ethical oversight and inadequate safety controls. Consider an AI system that suggests experimental protocols in biomedical research (e.g., chemical components with unknown toxicity) without the rigorous human review

⁷<https://research.google/blog/accelerating-scientific-breakthroughs-with-an-ai-co-scientist/>

needed to identify potential risks. This could lead to experiments that pose unforeseen dangers or violate established ethical standards.

3.2.4 Domains of Application. Regarding domains of interest, previous works have addressed idea and hypothesis in NLP, Engineering, Physics, Chemistry, Social Science and Medicine. Similarly, automated experimentation also relies on domain-specific datasets to guide the process of designing and testing experiments. In contrast, methods are typically domain-agnostic, though they are often developed and evaluated in specific domains. These methods address fundamental issues in LLMs, such as long-context inputs, post-tuning strategies, and model hallucination, making them potentially applicable across multiple domains.

3.2.5 Limitations and Future Directions. A large-scale study [206] comparing human researchers and LLMs finds that LLMs generate ideas judged to be more novel but slightly less feasible, highlighting challenges like limited diversity and self-evaluation failures. Additionally, given that ideas and hypotheses are theoretical and costly to validate, it is unclear whether they could lead to scientific discovery. Furthermore, previous methods lack due diligence through data, and therefore generated ideas and hypotheses are often too general [253]. Moreover, LLMs may generate recently discovered ideas and hypotheses, as they lack access to recent scientific papers [148]. Their outputs are very sensitive to the framing of input prompts [170]. Future work should focus on improving feasibility and diversity of ideas and hypotheses, incorporating real-time scientific papers, refining ideation and hypothesis generation through data inspection.

Automated experimentation with LLMs faces several additional challenges. First, LLMs often make critical errors, such as hallucinating results or outputting invalid references, which disrupt the precise steps required for experimental workflows. Another significant limitation is that LLMs struggle to integrate and align different modalities, such as video, audio, or sensory data, which are increasingly essential in modern scientific experimentation. Moreover, LLMs lack the critical analysis capabilities necessary to identify flaws or refine hypotheses during experimentation. In highly specialized scientific domains such as biology and chemistry, they may also struggle with precise reasoning and tool usage, which are vital for ensuring experimental success [187].

3.3 Text-based Content Generation

Under text-based content generation for science, we subsume different tasks generating specific text-based subparts of a scientific paper, such as automatically generating (i) the title, (ii) the abstract, (iii) the related work section, as well as (iv) citation generation. Also, frameworks aiming to automate the full paper writing process will be discussed, as well as using AI systems for subtasks such as proof-reading, paraphrasing, and press release generation.

3.3.1 Data. Open access research articles are a valuable data source for text-based content generation. These include scientific publisher repositories offering at least some open access content (e.g., [Nature portfolio](#), [Taylor & Francis](#)) as well as preprint repositories (e.g., [arXiv](#), [bioRxiv](#)). These open access repositories can be leveraged to develop datasets with pairs of titles and abstracts or abstract and conclusion/future work pairs. Wang et al. [238] for example extract (i) title to abstract pairs, (ii) abstract to conclusion and future work pairs, and (iii) conclusion and future work to title pairs from PubMed. Annotated, task-specific datasets for scientific text generation are presented in Table 3.

3.3.2 Methods and Results. In the following, we survey approaches to generating textual content for science, such as title, abstract, related work and bibliography. An overview of these processes is given in Appendix A.2.3.

Dataset	Size	Sources	Application
Abstract-title humor annotated dataset [34]	2,638 humor annotated titles	ML & NLP domain	Title generation
PaperRobot [238]	27,001 title-abstract pairs; 27,001 abstract-conclusion & future work pairs; 20,092 conclusion & future work-title pairs	PubMed	Title generation, abstract generation, conclusion & future work generation
ScisummNet [254]	1,000 papers + 20 citation sentences each	ACL Anthology Network	Related work generation
CORWA [136]	927 related work sections	NLP domain	Related work generation
CiteBench [69]	358,765 documents + citations	multiple, e.g., arXiv.org	Related work generation
SciTechNews [24]	2,431 papers + press releases	ACM TechNews	Press release generation

Table 3. Annotated or task-specific datasets for scientific text generation.

Title Generation. Generating adequate titles for scientific papers is an important task because titles are the first access point of a paper and can attract substantial reader interest; titles can also influence the reception of a paper [125]. Consequently, several works have targeted generating titles automatically, often using paper abstracts as input. For example, Mishra et al. [160] use a pipeline of three modules, viz. generation by transformer based (GPT2) models, selection (from multiple candidates) and refinement. Chen and Eger [34] also leverage transformers for title generation from abstracts but they in addition allow for generation of humorous titles (which may be even more impactful) when an input flag is set appropriately. To achieve this, they annotate a training dataset of humorous titles from the fields of machine learning and NLP. They explore different models including BART, GPT2, and T5 besides the more recent ChatGPT-3.5 LLM, finding that none of them can adequately generate humorous titles. They also explore generating titles from full texts instead of abstracts, with mixed results. Wang et al. [238] address the problem differently by drafting title names based on future work sections of previous related papers.

Abstract Generation. There are several approaches trying to assess the capabilities of proprietary LLMs to generate abstracts based on context information such as paper titles, journal names, keywords or the full text of the paper. Hwang et al. [99] assess the ability of GPT 3.5 and GPT 4 to generate abstracts based on a full text. The results are manually evaluated using the Consolidated Standards of Reporting Trials for abstracts, a standard published with an aim to enhance the overall quality of scientific abstracts [90]. While the readability of abstracts generated by GPT is rated higher, their overall quality is inferior to the original abstracts. Also, minimal errors are reported in the AI generated abstracts. Wang et al. [238] generate abstracts from titles, leveraging transformers and knowledge bases. Also generating abstracts from titles, Gao et al. [70] collect 50 research abstracts from five medical journals and apply ChatGPT to generate research abstracts based on their titles and the name of one of the five journals. The original and the generated abstracts are then evaluated with AI output detectors and with blinded human reviewers to identify which of the abstracts are automatically generated. Human reviewers are able to identify 68% of the generated abstracts as being automatically generated, but also incorrectly identify 14% of original abstracts as being LLM generated. Applying AI output detectors, most generated abstracts can be identified by the GPT-2 Output Detector assigning a median of 99.98% generated scores to generated abstracts and a median 0.02% to original abstracts. However, Anderson et al. [6] have shown that after automatically paraphrasing AI generated text, the performance of AI detectors such as GPT-2 Output Detector decrease drastically. Farhat et al. [63] evaluate the performance of ChatGPT generating abstracts based on 3 keywords, the name of a database (Scopus or web of science) and the task to analyze bibliographic data in the domain indicated by the keywords. The authors then compare the generated abstract to an actual abstract on the same

topic. After a manual comparison of the results, the authors come to the conclusion that at the time the study was conducted, ChatGPT is not a trustworthy tool for retrieving and assessing bibliographic data.

Long Text Generation. Some approaches aim at automating the full paper writing process. The **AI Scientist** [151] presents a comprehensive framework designed to support the entire scientific research cycle, encompassing tasks such as idea generation, hypothesis formulation, experimental planning, and execution. While its primary focus is not on long-form text generation, AI Scientist is able to generate entire scientific papers. By incorporating structured scientific knowledge (e.g. experimental results), the framework can draft papers that adhere to domain-specific requirements, involving the integration of relevant citations and conforming to disciplinary norms. Despite its ability to produce comprehensive paper drafts, the framework does not explicitly address the challenge of maintaining coherence across extended narratives, and their dependencies. **LongWriter** [9] and **LongEval** [245] directly address the challenge of generating extended text by introducing architectural modifications aimed at enhancing coherence and structural consistency in long-form outputs. The framework employs hierarchical attention mechanisms to ensure thematic consistency across long text and applies fine-tuning strategies to align outputs with user prompts. LongWriter conducts experiments on several domains, including academic and monograph texts. For academic content, the model can generate structured arguments and effectively incorporate domain-specific terminologies. However, noticeable issues remain around factual consistency, the integration of citations, and redundancy in the generated text. However, by conducting experiments on various models in academic, wikipedia and blog domains, LongEval shows that the larger models trained with general instruction data performs similar to those specifically trained (e.g., LongWriter). **LongReward** [262] leverages reinforcement learning to improve long-text generation. The model employs custom reward signals that prioritize coherence, factual accuracy, and linguistic quality. These reward mechanisms are particularly relevant for scientific text generation, where accuracy and adherence to domain-specific conventions are crucial.

Related Work Generation. Already in the past, there has been a substantial body of work on related work generation through text summarization, most of which differ in their approach (extractive or abstractive) and the length of citation text (sentence-level or paragraph-level). Extractive approaches focus on selecting sentences from cited papers and reordering the extracted sentences to form a paragraph of related work. For instance, Hoang and Kan [89] propose an extractive summarization approach that selects sentences describing the cited papers to generate the related work section of a target paper. This approach relies on the full text of the target paper. Subsequent extractive approaches differ from this approach in how they order the extracted sentences: While Wang et al. [242], Chen and Zhuge [33], and Wang et al. [237] assume that the sentence order is given, Hu and Wan [94] and Deng et al. [50] take advantage of an automatic approach to reorder sentences based on topic coherence. However, extractive approaches often struggle to produce coherent text, as they simply concatenate sentences without ensuring a cohesive narrative flow. In contrast, abstractive related work generation leverages devices of rewriting and restructuring to generate a summary of a cited paper. Most of the abstractive approaches are based on language models and focus on either generating (a) a single sentence from a single reference or (b) a paragraph from multiple references. Typically, the abstractive process is repeated multiple times until a related work section is complete. AbuRa'ed et al. [1] introduce an abstractive summarization approach to generate citation sentences in a single-reference setup. Their approach has been trained on the **ScisummNet** corpus with paper abstracts as inputs and citation sentences as outputs. Li et al. [136] further extend this idea to a multiple-reference setup, namely generating a paragraph of citation sentences from various cited papers. Their approach has been trained on the **CORWA** corpus to generate both citation and transition sentences. Additionally, instead of using paper abstracts as inputs, Li et al. [135] propose to retrieve relevant sentences from cited papers to generate citation sentences. More

recently, works such as Şahinuç et al. [193] explore instruction promoting with LLMs, which is alternative to extractive and abstractive approaches, to generate citation sentences. Overall, extractive approaches, while factual, often lack fluency and coherence. In contrast, abstractive approaches and instruction prompting, which are based on (large) language models, do not struggle with these issues, however, they suffer from factual errors, known as hallucination.

Citation Generation. Bibliographic references in scientific papers are important components for ensuring the scientific integrity of the authors. However, in many cases, cited articles of bibliographic references generated by LLMs such as ChatGPT are reported not to exist, that is, are hallucinated or incorrect [63, 96, 134, 137]. Most of the studies reporting hallucinated or erroneous bibliographic references are case studies presenting one or more examples. Walters and Wilder [233], however, present a study in which they use ChatGPT-3.5 and ChatGPT-4 to produce 84 documents (short reviews of the literature) on 42 multidisciplinary topics. The resulting documents contain 636 bibliographic citations, which are further analyzed for errors and hallucinations. Their results show that 55% of the GPT-3.5 citations but only 18% of the GPT-4 citations are fabricated. Of the actual existing (non-fabricated) GPT-3.5 citations, 43% include substantive citation errors, and of the non-fabricated GPT-4 citations it is 24%. Even though this is a major improvement from GPT-3.5 to GPT4, problems with fabrication and errors in bibliographic citations remain. Therefore, for generated citations and references, it is of particular importance to ensure their accuracy and completeness.

Proof-reading and Paraphrasing. LLMs such as ChatGPT have been reported to provide useful assistance for scientific writing with regards to proof-reading and language review in order to enhance the readability of the paper. Subtasks these models can provide support for during the writing process include providing suggestions for improving the writing style, or proof-reading [195]. Additionally, some authors emphasize that LLMs can be helpful especially for non-native English speakers with regards to grammar, sentence structure, vocabulary and even translation, i.e., providing an English editing service [25, 97, 110]. Most papers on this topic are case studies, illustrating their research questions with one or more examples and their results are qualitatively evaluated by a human expert (typically the author of the paper). Hassanipour et al. [84] evaluate the effectiveness of ChatGPT in rephrasing not for improving the writing style, but for reducing plagiarism in the process of scientific paper writing. The results showed that even with explicit instructions to paraphrase or reduce plagiarism, the plagiarism rate remained relatively high.

Press Release Generation. Several studies attempt to generate press release articles for the general public based on scientific papers. Cao et al. [23] construct a manually annotated dataset for expertise-style transfer in the medical domain and apply various style transfer and sentence simplification models to convert expert-level language into layman’s terms. Goldsack et al. [79] develop standard seq-to-seq models to generate news summaries for scientific articles. Lastly, Cardenas et al. [24] propose a framework that integrates metadata from scientific papers and scientific discourse structures to model journalists’ writing strategies.

3.3.3 Ethical Concerns. In scientific work, authorship and plagiarism in AI generated texts are major concerns. In general, it is a challenge to distinguish between AI generated and human generated texts. Although there is a number of tools to detect AI-generated text (e.g., GPTZero or Hive), Anderson et al. [6] show that after applying automatic paraphrasing to AI generated text, the probability of a text to be human generated, increases. Therefore it is not possible to reconstruct if a text is an original work from a scientist or has been generated by an AI. In addition, it is also found that ChatGPT generated texts easily pass plagiarism detectors [3, 61]. Moreover, Macdonald et al. [155] raise the concern that the frequent use of LLMs for drafting research articles might lead to similar paragraphs and structure of many

papers in the same field. This again raises the question whether there should be a threshold for the acceptable amount of AI-generated content in scientific work [155].

3.3.4 Domains of Application. Text-based content generation is relevant for all scientific domains. Liang et al. [141] conduct a large-scale analysis across 950,965 papers published between January 2020 and February 2024 to measure the prevalence of LLM modified content over time. The papers they investigated were published on (i) arXiv including the five areas Computer Science, Electrical Engineering and Systems Science, Statistics, Physics, and Mathematics, (ii) bioRxiv, and (iii) Nature portfolio. Their results show the largest and fastest growth in Computer Science with up to 17.5% of the papers containing LLM modified content and the least LLM modifications in Mathematics papers (up to 6.3%). However, according to the Natural Language Learning & Generation arXiv report from September 2024, top-cited papers show notably fewer markers of AI-generated content compared to random samples [124].

3.3.5 Limitations and Future Directions. Numerous studies have investigated text-based content generation for the scientific domain and have shown their potential to assist scientists in different phases of writing a paper. While for some tasks such as proof-reading and paraphrasing, its capabilities are well established, others pose limitations. Therefore it is crucial that automatically generated text is always assessed by a human expert. Factual consistency and truthfulness are issues which need to be reviewed by a human in the loop for all types of text-based generated content. Current proprietary LLMs for example struggle in particular with generating existing and correct bibliographic citations. However, LLMs are advancing rapidly and studies evaluating LLMs are quickly outdated. Still, several ethical issues arise when text-generating systems are included in the scientific writing process, such as authorship, plagiarism, bias, and truthfulness. Therefore, in future research a focus on trustworthy, ethical AI systems is required.

3.4 Multimodal Content Generation and Understanding

Multimodal content generation in the scientific domain refers to generating multimodal scientific content such as figures and tables in scientific papers or, e.g., slides and posters in a post-publication process. Automatizing such tasks via AI is important for multiple reasons: (i) generating high-quality figures, tables, slides and posters is *difficult and time-consuming* (dimension: cost); (ii) high-quality multimodal content in a paper can have a large **effect** (dimension: benefit for authors) on citation or acceptance decisions [123]; (iii) tables, figures, posters and slides make scientific content easier accessible for a scientific audience and often represent compact representations of research results (dimension: benefit for readers). **Multimodal scientific content understanding** refers to understanding scientific images and tables, e.g., answering questions about the multimodal scientific content, providing captions or summaries for scientific figures and tables. Automatizing this understanding process promises to allow to automatically describe such multimodal objects, which can likewise be time-consuming and costly for human authors, and helps readers digest the content more easily.

3.4.1 Data. Table 4 provides an overview over datasets for multimodal content generation and understanding. Below, we give details on individual datasets and benchmarks.

Scientific Figure Understanding. Scientific figure understanding benchmarks typically contain QA pairs for scientific figures or summaries of figures. Kembhavi et al. [108] provide a dataset with over 5k richly annotated diagrams and over 15k questions and answers. Kahou et al. [103] introduce **FigureQA**, a synthetic dataset of over 100k scientific-style figures from five classes: (dot-)line plots, vertical and horizontal bar graphs, and pie charts. Associated with the images

Dataset	Size	Data Sources
Scientific Figure Understanding		
ArxivCap [130] FigureQA [103] ChartQA [159] CharXiv [243] ArxivQA [130] SPIQA [177] ChartSumm [248] SciMMIR[246]	6.4m images and 3.9m captions from 572k papers > 100k scientific-style figures 4.8k charts, 9.6k QA pairs 2.3k charts with descriptive and reasoning questions 35k figures with 100k QA pairs 152k figures with 270k QAs 84k charts 530K figures and tables image-text pairs	Arxiv Synthetic statista.com, pewresearch.com, etc. Arxiv Arxiv 19 top-tier academic conferences Knoema Arxiv
Scientific Figure Generation		
DaTikZ [14, 15] SciImage [263] SciDoc2DiagramBench [161] ChartMimic [204]	180k-360k pairs of text captions and Tikz code 404 instructions and 3k generated scientific images 1,080 Extrapolated-Diagrams with the format of "<paper(s), intent of diagram, gold diagram>" 1000 triplets of (figure, instruction, code) instances	Mostly Arxiv and TeX.StackExchange Manual (template) construction ACL Anthology Physics, Computer Science, Economics, etc.
Scientific Table Understanding		
SciGen [163] NumericNLG [216] SciXGen [32]	1.3k pairs of scientific tables and their descriptions 1.3k pairs of scientific tables and their descriptions 484k tables from 205k papers	Arxiv (especially cs.CL and cs.LG) ACL Anthology Arxiv
Scientific Table Generation		
ArXivDigestTables [167]	2,228 literature review tables extracted from arXiv papers that synthesize a total of 7,542 research paper	literature review tables from ArXiv papers from April 2007 until November 2023
Scientific Slides and Poster Generation		
SciDuet [217] DOC2PPT [68] Persona-Aware-D2S [162]	1,088 papers and 10,034 slides by their authors 5,873 papers and 98,856 slides by their authors 75 papers from SciDuet, and 300 slides	NeurIPS/ICML/ACL Anthology CV (CVPR, ECCV, BMVC), NLP (ACL, NAACL, EMNLP), ML (ICML, NeurIPS, ICLR) ACL Anthology

Table 4. Multimodal content generation and understanding datasets.

are more than 1m questions generated using 15 different templates. Later research focuses on harder and more realistic QA pairs. Masry et al. [159] present **ChartQA**, which provides complex reasoning questions over charts sourced from various sources related to science. Wang et al. [243] introduce **CharXiv**, a manually curated dataset of 2.3k "natural, challenging, and diverse" charts from Arxiv papers and both descriptive and reasoning questions for them. Li et al. [130] introduce **ArxivQA**, a dataset of 35k scientific figures sourced from Arxiv for which GPT4o generates 100k QA pairs after manual filtering. Pramanick et al. [177] present **SPIQA**, a dataset of 270k manually and automatically created QA pairs to interpret complex scientific figures and tables. In contrast to focusing on question-answering for scientific figures, Xu et al. [248] consider the chart summarization problem and a dataset comprising more than 190k instances building on top of existing datasets such as **ChartSumm** [185], which contains more than 84k charts along with their metadata and descriptions covering a wide range of topics and chart types to generate short and long summaries.

Scientific Figure Generation. Recently, several datasets for scientific figure generation have been proposed. Belouadi et al. [14, 15] propose **DaTikZ** and **DaTikZ-v2** which contain (augmented) captions of scientific figures as instructions along with corresponding TikZ code, sourced from Arxiv submissions. Belouadi et al. [15] also provide **SketchFig**, a benchmark of 549 figure-sketch pairs to convert sketches into scientific figures; SketchFig is sourced from TEX stackexchange. Shi et al. [204] provide **ChartMimic**, a manually curated benchmark of 1000 triplets of (instruction, code, figure) instances for Chart generation across various domains (Physics, Computer Science, Economics, etc.).

ChartMimic is obtained by having human annotators write Python code for prototype charts. Zhang et al. [263] provide **ScImage**, which contains targeted template instructions focusing on different understanding dimensions (spatial, numeric, attribute). For a subset of the data, the authors also provide reference figures, evaluated as being of high-quality by human annotators. In contrast to the other benchmarks, ScImage also contains instructions in non-English languages. Mondal et al. [161] provide **SciDoc2DiagramBench**, a benchmark comprising 1,000 extrapolated diagrams paired with 89 ACL papers, along with human-written intents. All diagrams are extracted from the corresponding presentation slides of these papers. The intents describe how the content from each paper can be translated into the extrapolated diagrams for presentation purposes. Luo et al. [152] provide **nvBench**, a benchmark of 25k tuples of natural language queries and corresponding visualizations. nvBench is based on 153 databases and contains more than 7k visualizations on seven chart types. nvBench is synthesized from natural language to SQL benchmarks.

Scientific Table Understanding. Table understanding often comes as table-to-text generation, which focuses on producing accurate textual descriptions that reflect table content. **SciGen** [163] and **numericNLG** [216] are benchmarks specifically focused on scientific table reasoning, both emphasizing arithmetic reasoning over numerical tables, containing 1.3k expert-annotated tables. The annotations include the tables and parts of the scientific papers that describe the corresponding findings of the annotated tables. A specific subtask of these benchmarks is explored in Ampomah et al. [5], which focuses on generating textual explanations for tables reporting ML model performance metrics. This dataset pairs numerical tables of classification performance (e.g., precision, recall, and accuracy) with expert-written textual explanations that analyze and interpret the metrics. Datasets like **HiTab** [38] tackle the complexity of hierarchical tables commonly found in statistical reports, introducing numerical reasoning tasks that require models to account for implicit relationships and hierarchical indexing within tables. **SciXGen** [32] broadens the scope of table-to-text generation with context-aware scientific text generation. By drawing from over 200k scientific papers, SciXGen requires models to generate descriptions for tables, figures, and algorithms, grounded in the surrounding body text.

Scientific Table Generation. Table generation often comes in the form of text-to-table generation [49, 101, 205], the process of converting unstructured textual information into structured tabular formats. This process is particularly valuable for scientific domains where textual data often contains detailed experimental results, observations, or findings that need transformation into structured tables. In the scientific domain, **ArxivDIGESTables** [167] addresses the specific challenge of automating the creation of literature review tables. Rows in these tables represent individual papers, while columns capture comparative aspects such as methods, datasets, and results. ArxivDIGESTables supports the generation of literature review tables by leveraging additional grounding context, such as captions and in-text references.

Scientific Slide and Poster Generation. Most early efforts to automatically generate presentation slides from scientific papers relied on relatively small datasets for system development and evaluation. For example, Sravanthi et al. [211] collect source code (.tex files and figures) from eight papers and generate presentations based on them. Similarly, Hu and Wan [93] and Wang et al. [240] utilize 1,200 paper-slide pairs and 175 paper-slide pairs, resp., within the CS domain. For scientific poster generation, Qiang et al. [183] construct a dataset of 25 pairs of scientific papers and their corresponding posters. However, these datasets are often inaccessible to the public due to various constraints.

Two open-source datasets have emerged as widely used resources for follow-up research for scientific slide generation: **DOC2PPT** [68] and **SciDuet** [217]. DOC2PPT [68] comprises 5,873 paired scientific documents and their associated presentation slide decks with around 100,000 slides, drawn from three research communities: computer vision (CVPR,

ECCV, BMVC), natural language processing (ACL, NAACL, EMNLP), and machine learning (ICML, NeurIPS, ICLR). SciDuet [217] comprises 1,088 papers and 10,034 slides from conferences such as ICML, NeurIPS, and the ACL Anthology.

3.4.2 Methods and Results. In the following, we survey approaches to multimodal content generation and understanding. A summary table, along with additional related works, is provided in Appendix A.2.4.

Scientific figure understanding. Scientific figure understanding is typically framed in terms of (visual) QA, e.g., whether models are able to adequately answer questions on a given input figure [108]. Several recent works train baseline models such as transformers [159] or alternatives [103] such as [196] as well as explore recent LLMs on benchmarks [130, 243]. They generally show a big gap between proprietary models like GPT4o and the strongest open-source models and a big gap of all models to human performance. For chart summarization, Rahman et al. [185] find that older PLMs such as BART and T5 suffer from hallucination and missing out of important data points. Xu et al. [248] propose ChartAdapter, a lightweight transformer module that can be combined with LLMs for improved modeling of chart summarization. **Evaluation** of scientific figure understanding benchmarks mostly leverages automatic metrics. For example, Xu et al. [248] report out-dated and unreliable metrics such as BLEU and ROUGE for evaluating chart summaries; Pramanick et al. [177] report both human and automatic evaluation, using traditional QA metrics such as Meteor, Rouge, and BERTScore and novel LLM based metrics.

Scientific Figure Generation. Early work in the context of visualization for science (and beyond) dates back to the 1980s and 1990s at least [156, 191, 192]. Subsequent research [165, 209, 220] further developed rule-based or hybrid approaches, involving parsers and grammars, while custom neural architectures were later also explored [51, 145, 224].

Most recent work leverages multimodal LLMs. While Maddigan and Susnjak [157] explore diverse pre-trained LLMs, such as ChatGPT and GPT3, Voigt et al. [229] investigate smaller LLMs for real-time graphics generation on a CPU. Belouadi et al. [14, 15] treat the problem as a TikZ code generation problem where the input is (i) a scientific caption [14] or (ii) a sketch or image [15]. Both works fine-tune custom LLMs on datasets leveraged from Arxiv. Shi et al. [204] aim to generate Python code from instructions and/or images, specifically focusing on charts. They evaluate 3 proprietary and 11 open-weight LLMs on their ChartMimic benchmark, finding that even the best models (GPT-4 and Claude-3-opus) have substantial room for improvement. Zhang et al. [263] provide a template based approach to evaluate various multimodal LLMs in generating scientific images. They explore LLMs that can generate code (TikZ and Python) and ones that directly generate images, without intermediate code synthesis and in addition consider different input languages (English, German, Chinese, Farsi). They find that, except for GPT4o, most models struggle substantially in generating adequate scientific images. Zala et al. [260] explore the diagram generation task where LLMs first generate diagram plans and then the diagrams themselves. Mondal et al. [161] explore the same task with an additional refinement (feedback from multiple critic models) to enhance factual correctness. **Evaluation** of the models comprises automatic metrics including DreamSim [67], for image similarity, crystal Bleu [60] for code similarity and ClipScore [87] for text-image similarity, and human evaluation by ‘domain experts’. The former are typically reported to have low or medium correlation with the latter, establishing the need for domain specific evaluation in future work.

Scientific Table Understanding. Table-to-text generation encompasses a range of methodologies designed to transform structured tabular data into coherent and accurate textual descriptions. These techniques process, reason over, and utilize tabular structures to address challenges such as logical reasoning, content fidelity, and domain-specific adaptation. *Serialization* is a foundational approach where tables are linearized into sequences compatible with transformer-based language models. In this method, tables are converted into linear text sequences using special characters to delineate

structure [7, 163, 169]. *Structure-aware methods* explicitly model the inherent relationships and hierarchies within tables to enhance reasoning and generation fidelity. These include (a) *Intermediate Representations* [133, 267, 268]; (b) *Structure-Aware Pretraining* [116, 172, 251]; (c) *Structure-Aware Self-Attention Mechanisms* [146, 234]. **Evaluation:** Common metrics like BLEU and BARTScore are widely used to evaluate the fluency and relevance of generated text against reference outputs. However, ensuring faithfulness to the source table remains a significant challenge, often requiring human evaluation for accurate assessment [163, 172].

Scientific Table Generation. While none of the existing approaches have yet been applied specifically to scientific table generation, several methodologies present promising directions. The gTBLS (Generative Tables) approach [221] proposes a two-stage table generation process. The first stage infers the table structure from input text, while the second stage generates table content by formulating table-guided questions; this enhances syntactic validity and logical coherence of generated tables. In the context of open-structure table extraction, OpenTE [53] tackles the task of extracting tables with intrinsic semantic, calculational, and hierarchical structure from unstructured text. OpenTE introduces a three-step pipeline that identifies semantic and relational connections among table columns, extracts structured data, and grounds the output by aligning extracted data with the source text and table structure. The **evaluation** of text-to-table generation should focus on structural accuracy, value fidelity, and semantic coherence. TabEval [186] provides a promising direction by introducing a decomposition-based framework that breaks tables into atomic statements and evaluates them using entailment-based measures, though comprehensive evaluation still requires further advancements.

Scientific Slide and Poster Generation. For scientific slide generation, early works typically relied on heuristic rule-based approaches [211]. Later, researchers began to leverage machine learning approaches to extract key phrases and their corresponding important sentences [93, 127, 240]. All the aforementioned works focus on extracting sentences or phrases from the given paper to serve as the slide text content (“extractive approach”). In contrast, Fu et al. [68] and Sun et al. [217] take a different approach by training sequence-to-sequence models to generate sentences for the slide text content (“abstractive approach”). With recent advancements, researchers have started utilizing (multimodal) LLMs for generating scientific presentation slides [11, 158, 162]. Notably, all existing approaches take an extractive approach, where they extract images or tables directly from the original papers rather than generating new ones [11, 68, 162, 211, 217]. Generating posters from scientific papers has received less attention. Previous work has mainly explored different machine learning-based methods for generating key content and panel layouts from data [183, 250]. **Evaluation:** For scientific slide generation, most works evaluate the effectiveness of proposed approaches using automatic evaluation metrics and conduct human evaluation. The most common automatic evaluation metric is ROUGE. Fu et al. [68] introduce the Longest Common Figure Subsequence, which measures the quality of figures in the generated slides; Text-Figure Relevance (TFR), which assesses the similarity between the text of the ground truth slide and the generated slide containing the same figure; and Mean Intersection over Union, which evaluates layout quality. Recent studies have also begun utilizing LLMs to assess the quality of the generated slides [11, 158]. For scientific poster generation, in addition to conducting user studies, Qiang et al. [183] also measure the mean-square error (MSE) of the panel parameters (e.g., panel size, aspect ratio).

3.4.3 Ethical Concerns. Ethical concerns relating to models for figure, table, slide and poster generation especially include that these tools are technically limited (e.g., they may hallucinate content, be factually incorrect, and not correspond to the authors’ intentions), which could be overlooked, ignored or maliciously abused by human authors.

Such tools could also be misused to attack the scientific process, by purposefully producing incorrect results (e.g., as a testcase for adequate reviewing). Such risks may also be relevant in an educational context, e.g., when students use such tools for preparing term papers or theses.

3.4.4 Domains of Application. Many recent datasets for multimodal content generation and understanding are from ArXiv and more generally the STEM domain (science, technology, engineering and mathematics). Models such as DeTikZify and AutomaTikZ have also been fine-tuned on such data. This indicates a limitation both in terms of application scenarios and model assessments, as these may perform worse when applied in cross-domain scenarios.

3.4.5 Limitations and Future Directions. Common limitations among all multimodal generating and understanding approaches discussed include: (1) comparatively small size of datasets for fine-tuning models; (2) models often perform considerably below human performance on recently proposed benchmarks; (3) this concerns especially open-source non-proprietary models; (4) models and benchmarks are often limited to STEM domains and particularly Arxiv; (5) models lack reasoning abilities; (6) evaluation of models is difficult, especially for generation models, and current automatic metrics are often unsatisfactory. Specific problems occur in subfields: for example, for table generation, input text may be very long, which constitutes a problem for many current LLMs; for slide generation, there are no approaches that can generate slides from multiple documents (e.g., for tutorials) or that generate content beyond a paper (which may be necessary to include relevant background); for figure generation, models like AutomaTikZ are trained on captions, which are often not appropriate for generating the corresponding figure (e.g., a caption may simply be ‘Proof of Theorem X’), leading to mismatch between input descriptions and output figures.

3.5 Peer Review

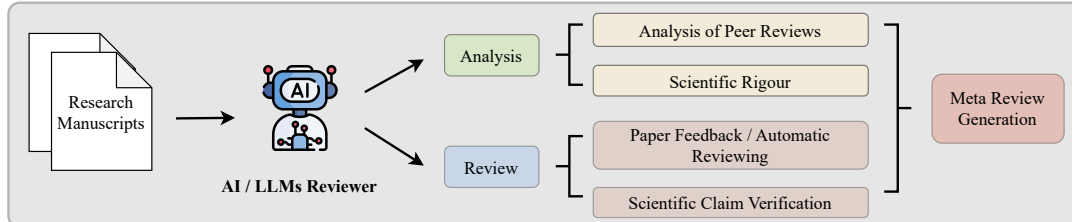


Fig. 3. Process of AI-enhanced peer review. In the analysis step, the LLM reviewer examines research manuscripts and evaluates peer reviews to assess scientific rigor. The review step involves providing feedback on the paper and verifying scientific claims. Finally, the gathered information is synthesized to generate a final meta-review.

The highest standard in scientific quality control is *peer reviewing*. In this process, the authors present their scientific argument (e.g., the findings of a study, a grant proposal, etc.), in form of a manuscript to their peers, who then assess its scientific validity and excellence. Often, this process has multiple stages, as shown in Fig. 3. For instance, in the ACL Rolling Review system,⁸ *reviewers* write detailed assessments. Afterwards, the *authors* may rebut the reviewers’ arguments and clarify questions to convince them to raise their scores. Finally, a meta-reviewer re-evaluates the whole scientific discussion and gives a final acceptance recommendation (which the overall program chairs may or may not adhere to). During this process, multiple (potentially multi-modal) artifacts are involved and created, mainly the

⁸<https://aclrollingreview.org>

manuscript under review, the written reviews, the author-reviewer discussion texts, and the meta-review. In general, peer review is considered a challenging, and subjective process, where reviewers are prone to unfair biases like sexism and racism, often relying on quick, simple heuristics [e.g., 188, 214]. At the same time, we are faced with an exploding number of submissions in some fields like AI [119], pushing peer reviewing systems to the limits of their capacities.

To counteract this problematic situation, researchers have worked on several problems under the umbrella of AI-supported peer review. Related overviews on the topic (or on some of its aspects) are given by [31, 54, 117, 120, 142, 213], pointing to the high relevancy of this problem. Here, we focus on existing works targeting the most established tasks, following the same structure as before, and provide an update on the recently published literature.

3.5.1 Data. Peer reviewing data is generally scarce, given that the scientific communities do not always make all reviewing artifacts publicly available under openly accessible licenses, with some exceptions like ICLR. As some exceptions, the PeerRead [104] collection of data from various sources (e.g., ACL, ICRL) and CiteTracked [174] are published along with citation information. As a prime example of how larger-scale open publishing of raw peer reviewing data may work, Dycke et al. [57] recently published the NLPeer corpus based on ARR reviews, for which they explicitly obtained the consent of the respective actors involved. For several tasks around peer review analyses,

Dataset	Size	Sources	Application
HedgePeer [74]	2,966 documents	ICLR 2018 reviews	Uncertainty detection
PolitePeer [17]	2,500 sentences	Various, e.g., ICLR	Politeness Analysis
COMPARE [207]	1,800 sentences	ICLR	Comparison Analysis
ReAct [40]	1,250 comments	ICLR	Actionability Analysis
MReD [203]	7,089 meta-reviews	ICLR	Meta-review analysis and generation
CiteTracked [174]	3,427 papers and 12k reviews	NeurIPS	citation prediction
MOPRD [143]	6,578 papers	PeerJ	Review Comment Generation
Revise and Resubmit [121]	5.4k papers	F1000Research	Tagging, Linking, Version Alignment
ORB [222]	92,879 reviews	OpenReview, SciPost	Acceptance Prediction
ARIES [48]	3.9k comments	OpenReview	Feedback-Edits Alignment, Revision Generation
DISAPERE [109]	506 review-rebuttal pairs	ICLR	review action analysis, polarity prediction, review aspect
PeerReviewAnalyze [73]	1,199 reviews	ICLR	Review Paper Section Correspondence, Paper Aspect Category Detection, Review Statement Role Prediction, Review Statement Significance Detection, and Meta-Review Generation
JitsuPeer [180]	9,946 review and 11,103 rebuttal sentences	ICLR	Argumentation Analysis, Canonical Rebuttal Scoring, Review Description Generation, End2End Canonical Rebuttal Generation

Table 5. Annotated or task-specific datasets for analyzing peer reviewing.

researchers have created annotated datasets. An overview of annotated and/ or task-specific datasets focusing on diverse aspects of peer review is provided in Table 5. Most recently, researchers focused on curating resources for supporting more complex tasks, like understanding the effect of peer review feedback on revisions of the manuscript [48] or on identifying the underlying attitudes that cause a specific criticism in peer review [180].

3.5.2 Methods and Results. Initial works were mostly based on more traditional machine learning methods and targeted simpler analyses involving sentence classification tasks. Later, deep learning approaches (also based on pre-trained

language models) and more complex analyses, e.g., argumentation analyses, were defining the state-of-the-art in computational peer review processing. Nowadays, researchers started exploring LLMs in prompting-based frameworks for complex tasks like peer review generation and meta-review generation.

Analysis of Peer Reviews. Prior works have analyzed peer reviews for a multitude of aspects, like uncertainty [74], politeness [17], and sentiment [27]. However, given that science as a whole and especially peer review relies to a large extent on convincing peers, large efforts have been spent on understanding arguments or argument-related aspects (e.g., substantiation of arguments) in peer review artifacts [e.g., 66, 95]. Here, most approaches leveraged pre-trained language models. For instance, Hua et al. [95] work on mining the arguments in peer reviews using conditional random fields, LSTMs, and BERT. In contrast, Guo et al. [81] and Fromm et al. [66] fully rely on (domain adjusted) pre-trained language models for argument mining like SciBERT, ArgBERT, and PeerBERT. Cheng et al. [37] leverage multi-task learning approaches based on LSTMs and BERT. In a similar vein, Purkayastha et al. [180] study the generation of rebuttals for author-reviewer discussions based on Jiu-Jitsu argumentation, a specific theory in argumentation theory.

Paper Feedback and Automatic Reviewing. Several works have explored methods to provide general feedback on scientific publications to fully or partially automate peer reviews. For instance, Li et al. [129] propose a multi-task learning approach for peer review score prediction, where different aspect score prediction tasks (e.g., novelty) can inform each other. Ghosal et al. [75] leverage the concept of sentiment to predict scores based on review texts. In a similar vein, Bharti et al. [18] leverage paper-review interactions to predict final decisions of a review process. Wang et al. [239] focus on explainability during review score prediction for several review categories by constructing knowledge graphs (e.g., one which represents the background of a paper). More recent works have included the generation of feedback texts into the problem setup. Bartoli and Medvet [13] frame the problem as exploring the potential of GPT-2 for conducting academic fraud by generating fake reviews. In contrast, Yuan et al. [259] ask whether it would be possible to automate reviewing leveraging targeted summarization models, a recently trending topic. For instance, Liu and Shah [149] explore prompting-based review generation with several LLMs like GPT-4, Vicuna, Llama. They find that GPT-4 performs best among the models tested and that task granularity matters. Similarly, Robertson [190] find GPT-4 to be “slightly” helpful for peer reviewing, and Liang et al. [140] demonstrate in a comparative study that users of a GPT-4-based peer review system found the feedback to be (very) helpful in more than half of the cases. D’Arcy et al. [47] show a multi-agent approach with LLMs that engage in a discussion to produce better results than a single model.

Scientific Rigor. Several attempts have been made to computationally analyze the rigor of scientific papers. For example, Soliman and Siponen [208] investigate how researchers use the word “rigor” in information system literature but discovered that the exact meaning was ambiguous in current research. Additionally, various automated tools have been proposed to assess the rigor of academic papers. Phillips [173] develop an online software that spots genetic errors in cancer papers, while Sun et al. [218] use knowledge graphs to assess the credibility of papers based on meta-data such as publication venue, affiliation, and citations. However, these methods are neither domain-specific nor do they provide sufficient guidance for authors to improve their narrative and writing. In contrast, SciScore [202] is an online system that uses language models to produce rigor reports for paper drafts, helping authors identify weaknesses in their presentation. More recently, James et al. [100] propose a bottom-up, data-driven framework that automates the identification and definition of rigor criteria while assessing their relevance in scientific texts. Their framework integrates three key components: rigor keyword extraction, detailed definition generation, and the identification of salient criteria. Additionally, its domain-agnostic design allows for flexible adaptation across different fields.

Scientific Claim Verification. The increasing number of publications requires the development of automated methods for verifying the validity and reliability of research claims. Scientific fact verification, which aims to assess the accuracy of scientific statements, often relies on external knowledge to support or refute claims [52, 228]. Several datasets have been developed to address this including SciFact-Open [232], which provides scientific claims and supporting evidence from abstracts. However, they are limited to the use of abstracts as the primary source of evidence. As the statements in abstract can also be inaccurate (e.g. overstated claims), it is important to evaluate the evidence in the main body of the paper to determine if the statements made in the abstract are well supported. On the other side, Glockner et al. [77, 78] propose a theoretical argumentation model to reconstruct fallacious reasoning of false claims that misrepresent scientific publications. The need to contextualize claims with supporting evidence is highlighted by Chan et al. [30], who introduce a dataset of claims extracted from lab notes. Unlike other datasets, this resource contains claims “actually in use”, providing a more realistic understanding of how researchers interact with scientific findings. The authors annotate these claims with links to figures, tables, and methodological details, and develop associated tasks to improve retrieval. While this provides valuable resources for context-based verification, it primarily focuses on factual verification and does not evaluate the potential for overstated claims. Beyond factual correctness, there is a growing recognition for the need to analyze how researchers present their findings. This includes the detection of overstatements, where authors exaggerate their achievements, and understatements, where the true impact of the research is downplayed [106]. Such analysis goes beyond the simple fact of a claim and is necessary to understand the presentation of a claim. Schlichtkrull et al. [200] present a qualitative analysis of how intended uses of fact verification are described in highly-cited NLP papers, particularly focusing on the introductions of the papers, to understand how these elements are framed. The work suggests that claims should be supported by relevant prior work and empirical results.

Meta Review Generation. Kumar et al. [118] tackle meta-review generation using a multi-encoder transformer network, and Li et al. [132] use a multi-task learning approach for refining pre-trained language models for the task. Stappen et al. [212] explore the aggregation of reviews for providing additional computational decision support to editors based on uncertainty-aware methods like soft labeling. Both Zeng et al. [261] and Santu et al. [197] rely on LLMs which they specifically prompt for the task.

3.5.3 Ethical Concerns. Given the critical role of scientific peer review for science, and, accordingly, for society as a whole, ethical considerations around AI-supported peer review are of utmost importance. As the general concerns around unfair biases in AI and the resulting harms apply [120], research on safe peer-reviewing support needs to be prioritized. For instance, von Wedel et al. [230] recently showed that LLMs exhibit affiliation biases when reviewing abstracts. In this context, any AI-support for peer reviewing needs to be critically evaluated [198], and solutions that target only a particular aspect in a collaborative environment that leaves the scientific autonomy to the human expert, may need to be preferred over end-to-end reviewing systems.

3.5.4 Domains of Application. Generally, peer review comes in many variations. For instance, the specific aspects to review for, how much textual content to produce, the specific scoring schemes, and the envisioned stages and dynamics of the reviewer and reviewer-author discussions may change. Thus, while none of the studies presented above targets a problem that is truly unique to any scientific domain, the particularities will likely be very different for each specific community and existing systems will need to be carefully evaluated before deployment.

3.5.5 Limitations and Future Directions. For existing studies on peer review, in particular, the variety of scientific domains that have been studied is still limited. As most of the works rely on data from OpenReview, most studies

focus on peer review within the ICLR and ACL communities [e.g., 40, 109]. To the best of our knowledge, for some domains, no single data set (yet, a data set further enriched with annotations or other additional information) exists (e.g., legal studies). Furthermore, scientific rigor, a critical aspect of peer review, remains underexplored. Most existing studies rely on predefined rigor checklists, such as those suggested by the NIH and MDAR [29], which are not easily scalable or transferable across different domains. Given these gaps, future research could benefit from exploring new domains of peer review, developing domain adaptation approaches, and advancing models for assessing scientific rigor. Additionally, in light of the ethical concerns discussed earlier, it is crucial to prioritize research on trustworthy AI support for peer review - ensuring that human experts retain autonomy in the review process.

4 Ethical Concerns

By now, there is a body of work addressing major ethical concerns related to generative AI. Baldassarre et al. [10], for instance, present a systematic literature review regarding the social impact of generative AI, especially taking into account 71 papers on ChatGPT. They identify the following areas of concern: privacy, inequality, bias, discrimination, and stereotypes. Another literature review on ethics and generative AI conducted by Hagendorff [82] identifies the following topics as becoming increasingly of interest: jailbreaking, hallucination, alignment, harmful content, copyright, models leaking private data, impacts on human creativity. The work also identifies 19 distinct clusters of ethics topics with fairness/bias being the most frequently mentioned, followed by safety, harmful content/toxicity, hallucinations, privacy, interaction risks, security/robustness on ranks two to six, and writing/research on rank 18. Ali and Aysan [2] review 364 recent papers on generative AI and ethics published from 2022 to 2024 in different domains including the use of generative AI in scientific research. Regarding academia, the prevalent topics identified as critical are the authenticity of the work, intellectual property and academic integrity. Sun et al. [219] argue that in application areas such as scientific research, ensuring the trustworthiness of LLMs is crucial. In particular truthfulness, i.e., the accurate representation of information, facts and results by an AI systems, is an essential challenge for LLMs.

Some benchmarks and datasets are designed to evaluate different aspects of truthfulness, for instance: TruthfulQA [144], HaluEval [128], and the FELM dataset [269] identify hallucinations. SelfAware [256] assesses the awareness of knowledge limitations. FreshQA [231] and Pinocchio [256] explore adaptability to rapidly evolving information. TrustLLM [219] is an extensive benchmark incorporating existing and new datasets on the six aspects truthfulness, safety, fairness, robustness, privacy, and machine ethics to assess the trustworthiness of LLMs. Their results show that, in general, proprietary LLMs (e.g., ChatGPT, GPT-4) outperform most open-source LLMs in trustworthiness — with Llama2 [225] as an exception. However, both proprietary LLMs, as well as Llama2 often struggled to provide truthful responses when relying solely on internal knowledge. However, their performance improved significantly with additional external knowledge. Moreover, the authors observed a positive correlation between trustworthiness and the functional effectiveness of the model in downstream tasks.

Editors of scientific publications are particularly challenged as the proportion of AI generated text in academic manuscripts is steadily increasing [36, 80, 114, 141] and AI models are also on the brink of being used in peer reviewing, cf. Section 3.5. The editors-in-chief of the Journal of Information Technology, for instance, elaborate in [199] on the limitations and risks of using generative AI in the production of scientific publications. They refer to the eight-point ‘Artificial Imperfections’ test to illustrate current limitations of generative AI: AI is (i) brittle, (ii) opaque, (iii) greedy, (iv) shallow and tone-deaf, (v) manipulative and hackable, (vi) biased, (vii) invasive, (viii) ‘faking it’. Nevertheless, the editors conclude that the use of AI should not be forbidden, however, if it is used, the authors must take full responsibility of the outcome, adhere to the “scientific principle of transparency” and give full and transparent disclosure of the usage

of AI in the respective publication, and moreover that “it is then up to the reviewers and editors to assess and make decisions on the specific use of that generative AI in a specific piece of research.” Similarly, Pu et al. [179] in their editorial to *iMeta* 3(2), 2024 state that AI-assisted technologies cannot be recognized as authors, the use of generative AI in a scientific manuscript/publication must be transparently disclosed, including the prompts, specific versions of the tools used. The authors are fully responsible for the integrity of their manuscripts. They must address ethical concerns and ensure the accuracy and fairness of AI-generated content, comply with data protection and privacy laws, and also consider copyright and intellectual property issues surrounding AI-generated content. The use of AI-generated images and multi-media must be specifically allowed. They explicitly prohibit the use of AI in the reviewing process.

5 Conclusion

In this paper, we surveyed approaches in the area of AI4Science, with a particular focus on recent large language model-based methods. We examined five key aspects of the research cycle: (1) search, (2) experimentation and research idea generation, (3) text-based content production, (4) multimodal content production, and (5) peer review. For each topic, we discussed relevant datasets, methods, and results, including evaluation strategies, while highlighting limitations and avenues for future research. Ethical concerns featured prominently in our survey, given the potential for misuse and challenges in maintaining scientific integrity in the face of AI-assisted content generation.

We hope that this survey inspires new initiatives in AI4Science, driving faster, more efficient, and more inclusive scientific discovery, experimentation, reporting and content synthesis—while upholding the highest ethical standards. As the ultimate goal of science is to serve humanity, we hope these advancements will accelerate knowledge creation and enhance the accessibility and reliability of research, leading to improved healthcare, medical treatments, economic processes, among a myriad of other societal benefits.

Acknowledgments

Yong Cao was supported by a VolkswagenStiftung Momentum grant. Jennifer D’Souza was supported by the [SCINEXT project](#) (BMBF, German Federal Ministry of Education and Research, Grant ID: 01IS22070). The NLLG Lab at UTN gratefully acknowledges support from the Federal Ministry of Education and Research (BMBF) via the research grant “Metrics4NLG” and the German Research Foundation (DFG) via the Heisenberg Grant EG 375/5-1. The work of Anne Lauscher is supported by the Excellence Strategy of the German Federal Government and the Federal States. Our AI use cases are document in the supplemental material.

References

- [1] Ahmed AbuRa’ed, Horacio Saggion, Alexander Shvets, and Àlex Bravo. 2020. Automatic related work section generation: experiments in scientific document abstracting. *Scientometrics* 125 (2020), 3159–3185.
- [2] Hassnian Ali and Ahmet Faruk Aysan. 2024. Ethical dimensions of generative AI: a cross-domain analysis using machine learning structural topic modeling. *Int. J. Ethics Syst.* (2024).
- [3] Signe Altmäe, Alberto Sola-Leyva, and Andres Salumets. 2023. Artificial intelligence in scientific writing: a friend or a foe? *Reproductive BioMedicine Online* 47, 1 (2023), 3–9.
- [4] Maha Amami, Gabriella Pasi, Fabio Stella, and Rim Faiz. 2016. An lda-based approach to scientific paper recommendation. In *Nat. Lang. Process. Inf. Syst. 21st Int. Conf. on Appl. Nat. Lang. to Inf. Syst. NLDB 2016, Salford, UK, June 22-24, 2016, Proc. 21*. Springer, 200–210.
- [5] Isaac Ampomah, James Burton, Amir Enshaei, and Noura Al Moubayed. 2022. Generating Textual Explanations for Machine Learning Models Performance: A Table-to-Text Task. In *Proc. Thirteen. Lang. Resour. Eval. Conf.*. European Language Resources Association, 3542–3551.
- [6] Nash Anderson, Daniel L Belavy, Stephen M Perle, Sharief Hendricks, Luiz Hespanhol, Evert Verhagen, and Aamir R Memon. 2023. AI did not write this manuscript, or did it? Can we trick the AI text detector into generated texts? The potential future of ChatGPT and AI in Sports & Exercise Medicine manuscript generation. *BMJ open sport & exercise medicine* 9, 1 (2023), e001568.

- [7] Ewa Andrejczuk, Julian Eisenschlos, Francesco Piccinno, Syrine Krichene, and Yasemin Altun. 2022. Table-To-Text generation and pre-training with TabT5. In *Find. Assoc. for Comput. Linguist. EMNLP 2022. ACL*, 6758–6766.
- [8] Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2024. Self-RAG: Learning to Retrieve, Generate, and Critique through Self-Reflection. In *The Twelfth Int. Conf. on Learn. Represent.*
- [9] Yushi Bai, Jiajie Zhang, Xin Lv, Linzhi Zheng, Siqi Zhu, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. 2024. Longwriter: Unleashing 10,000+ word generation from long context llms. *arXiv preprint arXiv:2408.07055* (2024).
- [10] Maria Teresa Baldassarre, Danilo Caivano, Berenice Fernandez Nieto, Domenico Gigante, and Azzurra Ragone. 2023. The social impact of generative ai: An analysis on chatgpt. In *Proc. 2023 ACM Conf. on Inf. Technol. for Soc. Good*. 363–373.
- [11] Sambaran Bandyopadhyay, Himanshu Maheshwari, Anandhavelu Natarajan, and Apoorv Saxena. 2024. Enhancing Presentation Slide Generation by LLMs with a Multi-Stage End-to-End Approach. In *Proc. 17th Int. Nat. Lang. Gener. Conf.*. ACL, 222–229.
- [12] Trapit Bansal, David Belanger, and Andrew McCallum. 2016. Ask the gru: Multi-task learning for deep text recommendations. In *proceedings 10th ACM Conf. on Recomm. Syst.*. 107–114.
- [13] Alberto Bartoli and Eric Medvet. 2020. *Exploring the Potential of GPT-2 for Generating Fake Reviews of Research Papers*.
- [14] Jonas Belouadi, Anne Lauscher, and Steffen Eger. 2024. AutomaTikZ: Text-Guided Synthesis of Scientific Vector Graphics with TikZ. In *Proc. ICLR*.
- [15] Jonas Belouadi, Simone Paolo Ponzetto, and Steffen Eger. 2024. DeTikZify: Synthesizing Graphics Programs for Scientific Figures and Sketches with TikZ. In *The Thirty-eighth Annu. Conf. on Neural Inf. Process. Syst.*
- [16] Chandra Bhagavatula, Sergey Feldman, Russell Power, and Waleed Ammar. 2018. Content-Based Citation Recommendation. In *Proc. 2018 Conf. North Am. Chapter Assoc. for Comput. Linguist. Hum. Lang. Technol. Vol. 1 (Long Pap., Marilyn Walker, Heng Ji, and Amanda Stent (Eds.))*. ACL, New Orleans, Louisiana, 238–251.
- [17] Prabhat Bharti, Meith Navlakha, Mayank Agarwal, and Asif Ekbal. 2023. PolitePEER: does peer review hurt? A dataset to gauge politeness intensity in the peer reviews. *Lang. Resour. Eval.* (05 2023), 1–23.
- [18] Prabhat Kumar Bharti, Shashi Ranjan, Tirthankar Ghosal, Mayank Agrawal, and Asif Ekbal. 2021. PEERAssist: Leveraging on Paper-Review Interactions to Predict Peer Review Decisions. In *Towards Open Trust. Digit. Soc.*. Springer International Publishing, 421–435.
- [19] Christine L. Borgman. 2007. *Scholarship in the Digital Age: Information, Infrastructure, and the Internet*. MIT Press.
- [20] Lutz Bornmann, Robin Haunschild, and Rüdiger Mutz. 2021. Growth Rates of Modern Science: a Latent Piecewise Growth Curve Approach to Model Publication Numbers from Established and New Literature Databases. *Humanit. Soc. Sci. Commun.* 8, 224 (2021).
- [21] Jennifer A. Byrne. 2016. Improving the Peer Review of Narrative Literature Reviews. *Res. Integr. Peer Rev.* 1, 12 (2016).
- [22] Fengyu Cai, Xinran Zhao, Tong Chen, Sihao Chen, Hongming Zhang, Iryna Gurevych, and Heinz Koepl. 2024. MixGR: Enhancing Retriever Generalization for Scientific Domain through Complementary Granularity. In *Proc. 2024 Conf. on Empir. Methods Nat. Lang. Process.*. ACL, 10369–10391.
- [23] Yixin Cao, Ruihao Shui, Liangming Pan, Min-Yen Kan, Zhiyuan Liu, and Tat-Seng Chua. 2020. Expertise Style Transfer: A New Task Towards Better Communication between Experts and Laymen. In *Proc. 58th Annu. Meet. Assoc. for Comput. Linguist.*. ACL, 1061–1071.
- [24] Ronald Cardenas, Bingsheng Yao, Dakuo Wang, and Yufang Hou. 2023. ‘Don’t Get Too Technical with Me’: A Discourse Structure-Based Framework for Automatic Science Journalism. In *Proc. 2023 Conf. on Empir. Methods Nat. Lang. Process.*. ACL, 1186–1202.
- [25] Andres Castellanos-Gomez. 2023. Good practices for scientific article writing with ChatGPT and other artificial intelligence language models. *Nanomanufacturing* 3, 2 (2023), 135–138.
- [26] Miaosen Chai, Emily Herron, Erick Cervantes, and Tirthankar Ghosal. 2024. Exploring Scientific Hypothesis Generation with Mamba. In *Proc. 1st Workshop on NLP for Sci. (NLP4Science)*. 197–207.
- [27] Souvic Chakraborty, Pawan Goyal, and Animesh Mukherjee. 2020. Aspect-based Sentiment Analysis of Scientific Reviews. In *Proc. ACM/IEEE Jt. Conf. on Digit. Libr. 2020 (Virtual Event, China) (JCDL ’20)*. Association for Computing Machinery, New York, NY, USA, 207–216.
- [28] Iain Chalmers, Larry V. Hedges, and Harris Cooper. 2002. A Brief History of Research Synthesis. *Eval. & Health Prof.* 25, 1 (2002), 12–37.
- [29] Karen Chambers, Andy Collings, Chris Graf, Veronique Kiermer, David Thomas Mellor, Malcolm Robert Macleod, Sowmya Swaminathan, Deborah Sweet, et al. 2019. Towards minimum reporting standards for life scientists. (2019).
- [30] Chu Sern Joel Chan, Aakanksha Naik, Matthew Akamatsu, Hanna Bekele, Erin Bransom, Ian Campbell, and Jenna Sparks. 2024. Overview of the Context24 Shared Task on Contextualizing Scientific Claims. In *Proc. Fourth Workshop on Sch. Document Process. (SDP 2024)*. 12–21.
- [31] A. Checco, L. Bracciale, P. Loreti, S. Pinfield, and G. Bianchi. 2021. AI-assisted peer review. *Humanit. Soc. Sci. Commun.* 8, 1 (January 2021), 1–11.
- [32] Hong Chen, Hiroya Takamura, and Hideki Nakayama. 2021. SciXGen: A Scientific Paper Dataset for Context-Aware Text Generation. In *Find. Assoc. for Comput. Linguist. EMNLP 2021. ACL*, 1483–1492.
- [33] Jingqiang Chen and Hai Zhuge. 2019. Automatic generation of related work through summarizing citations. *Concurr. Comput. Pract. Exp.* 31, 3 (2019), e4261.
- [34] Yanran Chen and Steffen Eger. 2023. Transformers Go for the LOLs: Generating (Humorous) Titles from Scientific Abstracts End-to-End. In *Proc. 4th Workshop on Eval. Comp. NLP Syst.*. ACL, 62–84.
- [35] Ziru Chen, Shijie Ning, Qianheng Zhang, Boshi Wang, Botao Yu, Yifei Li, Zeyi Liao, Chen Wei, Zitong Lu, Vishal Dey, Mingyi Xue, Frazier N. Baker, Benjamin Burns, Daniel Adu-Ampratwum, Xuhui Huang, Xia Ning, Song Gao, Yu Su, and Huan Sun. 2024. ScienceAgentBench: Toward Rigorous Assessment of Language Agents for Data-Driven Scientific Discovery. *arXiv:2410.05080* [cs.CL]

- [36] Huzi Cheng, Bin Sheng, Aaron Lee, Varun Chaudhary, Atanas G Atanasov, Nan Liu, Yue Qiu, Tien Yin Wong, Yih-Chung Tham, and Ying-Feng Zheng. 2024. Have AI-Generated Texts from LLM Infiltrated the Realm of Scientific Writing? A Large-Scale Analysis of Preprint Platforms. *bioRxiv* (2024), 2024–03.
- [37] Liying Cheng, Lidong Bing, Qian Yu, Wei Lu, and Luo Si. 2020. APE: Argument Pair Extraction from Peer Review and Rebuttal via Multi-task Learning. In *Proc. 2020 Conf. on Empir. Methods Nat. Lang. Process. (EMNLP)*. ACL, 7000–7011.
- [38] Zhoujun Cheng, Haoyu Dong, Zhiruo Wang, Ran Jia, Jiaqi Guo, Yan Gao, Shi Han, Jian-Guang Lou, and Dongmei Zhang. 2022. HiTab: A Hierarchical Table Dataset for Question Answering and Natural Language Generation. In *Proc. 60th Annu. Meet. Assoc. for Comput. Linguist.*, Vol. 1. ACL, 1094–1110.
- [39] Yizhou Chi, Yizhang Lin, Sirui Hong, Duyi Pan, Yaying Fei, Guanghao Mei, Bangbang Liu, Tianqi Pang, Jacky Kwok, Ceyao Zhang, et al. 2024. SELA: Tree-Search Enhanced LLM Agents for Automated Machine Learning. *arXiv preprint arXiv:2410.17238* (2024).
- [40] Gautam Choudhary, Natwar Modani, and Nitish Maurya. 2021. *ReAct: A Review Comment Dataset for Actionability (and more)*. Springer International Publishing, 336–343.
- [41] John Clement. 1989. Learning via Model Construction and Criticism: Protocol Evidence on Sources of Creativity in Science. In *Handbook of Creativity: Assessment, Theory and Research*, John A. Glover, Royce R. Ronning, and Cecil R. Reynolds (Eds.). Plenum, 341–381.
- [42] John J. Clement. 2022. Multiple Levels of Heuristic Reasoning Processes in Scientific Model Construction. *Front. Psychol.* 13 (2022).
- [43] Arman Cohan, Sergey Feldman, Iz Beltagy, Doug Downey, and Daniel Weld. 2020. SPECTER: Document-level Representation Learning using Citation-informed Transformers. In *Proc. 58th Annu. Meet. Assoc. for Comput. Linguist.*, Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (Eds.). ACL, Online, 2270–2282.
- [44] Alessandro Conti, Enrico Fini, Paolo Rota, Yiming Wang, Massimiliano Mancini, and Elisa Ricci. 2024. Automatic benchmarking of large multimodal models via iterative experiment programming. *arXiv preprint* (jun 2024).
- [45] Cristina Cornelio, Sanjeeb Dash, Vernon Austel, Tyler R. Josephson, Joao Goncalves, Kenneth L. Clarkson, Nimrod Megiddo, Bachir El Khadir, and Lior Horeish. 2023. Combining Data and Theory for Derivable Scientific Discovery with AI-Descartes. *Nat. Commun.* 14, 1777 (2023).
- [46] Paul Covington, Jay Adams, and Emre Sargin. 2016. Deep neural networks for youtube recommendations. In *Proc. 10th ACM conference on recommender systems*. 191–198.
- [47] Mike D’Arcy, Tom Hope, Larry Birnbaum, and Doug Downey. 2024. MARG: Multi-Agent Review Generation for Scientific Papers. *arXiv preprint arXiv:2401.04259* (2024).
- [48] Mike D’Arcy, Alexis Ross, Erin Bransom, Bailey Kuehl, Jonathan Bragg, Tom Hope, and Doug Downey. 2023. ARIES: A Corpus of Scientific Paper Edits Made in Response to Peer Reviews. *arXiv preprint arXiv:2306.12587* (2023).
- [49] Zheyue Deng, Chunkit Chan, Weiqi Wang, Yuxi Sun, Wei Fan, Tianshi Zheng, Yauwai Yim, and Yangqiu Song. 2024. Text-Tuple-Table: Towards Information Integration in Text-to-Table Generation via Global Tuple Extraction. In *Proc. 2024 Conf. on Empir. Methods Nat. Lang. Process.*. ACL, 9300–9322.
- [50] Zekun Deng, Zixin Zeng, Weiye Gu, Jiawen Ji, and Bolin Hua. 2021. Automatic Related Work Section Generation by Sentence Extraction and Reordering. (2021).
- [51] Victor Dibia and Çağatay Demiralp. 2019. Data2Vis: Automatic Generation of Data Visualizations Using Sequence-to-Sequence Recurrent Neural Networks. *IEEE Comput. Graph. Appl.* 39, 5 (2019), 33–46.
- [52] Alphaeus Dmonte, Roland Oruche, Marcos Zampieri, Prasad Calyam, and Isabelle Augenstein. 2024. Claim Verification in the Age of Large Language Models: A Survey. *arXiv preprint arXiv:2408.14317* (2024).
- [53] Haoyu Dong, Mengkang Hu, Qinyu Xu, Haochen Wang, and Yue Hu. 2024. OpenTE: Open-Structure Table Extraction From Text. In *ICASSP 2024 - 2024 IEEE Int. Conf. on Acoust. Speech Signal Process. (ICASSP)*. 10306–10310.
- [54] Iddo Drori and Dov Te’eni. 2024. Human-in-the-Loop AI Reviewing: Feasibility, Opportunities, and Risks. *J. Assoc. for Inf. Syst.* 25, 1 (2024), 98–109.
- [55] John A. Drozdz and Michael R. Ladomery. 2024. The Peer Review Process: Past, Present, and Future. *Br. J. Biomed. Sci.* 81 (2024).
- [56] Jennifer D’Souza, Sören Auer, and Ted Pedersen. 2021. SemEval-2021 Task 11: NLPContributionGraph - Structuring Scholarly NLP Contributions for a Research Knowledge Graph. In *Proc. 15th Int. Workshop on Semantic Eval. (SemEval-2021)*, Alexis Palmer, Nathan Schneider, Natalie Schluter, Guy Emerson, Aurelie Herbelot, and Xiaodan Zhu (Eds.). ACL, Online, 364–376.
- [57] Nils Dycke, Ilia Kuznetsov, and Iryna Gurevych. 2023. NLPeer: A Unified Resource for the Computational Study of Peer Review. In *Proc. 61st Annu. Meet. Assoc. for Comput. Linguist.*, Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (Eds.), Vol. 1. ACL, Toronto, Canada, 5049–5073.
- [58] Sašo Džeroski and Ljupčo Todorovski. 2007. *Computational Discovery of Scientific Knowledge: Introduction, Techniques, and Applications in Environmental and Life Sciences*. Number 4660. Springer.
- [59] D. N. Ege, M. Goudswaard, J. Gopsill, B. Hicks, and M. Steinert. 2023. The IDEA Challenge 2022 dataset. *Zenodo* (2023).
- [60] Aryaz Eghbali and Michael Pradel. 2023. CrystalBLEU: Precisely and Efficiently Measuring the Similarity of Code. In *Proc. 37th IEEE/ACM Int. Conf. on Autom. Softw. Eng.*. Association for Computing Machinery, Article 28, 12 pages.
- [61] Holly Else. 2023. Abstracts written by ChatGPT fool scientists. *Nature* 613 (2023), 423.
- [62] Thomas Elsken, Jan Hendrik Metzen, and Frank Hutter. 2019. Neural architecture search: A survey. *J. Mach. Learn. Res.* 20, 55 (2019), 1–21.
- [63] Faiza Farhat, Shahab Saquib Sohail, and Dag Øivind Madsen. 2023. How trustworthy is ChatGPT? The case of bibliometric analyses. *Cogent Eng.* 10, 1 (2023), 2222988.
- [64] Ronald A. Fisher. 1925. *Statistical Methods for Research Workers*. Oliver and Boyd.

- [65] Ronald A. Fisher. 1935. *The Design of Experiments*. Oliver and Boyd.
- [66] Michael Fromm, Evgeniy Faerman, Max Berrendorf, Siddharth Bhargava, Ruoxia Qi, Yao Zhang, Lukas Dennert, Sophia Selle, Yang Mao, and Thomas Seidl. 2021. Argument Mining Driven Analysis of Peer-Reviews. *Proc. AAAI Conf. on Artif. Intell.* 35, 6 (May 2021), 4758–4766.
- [67] Stephanie Fu, Netanel Y. Tamir, Shobhita Sundaram, Lucy Chai, Richard Zhang, Tali Dekel, and Phillip Isola. 2023. DreamSim: learning new dimensions of human visual similarity using synthetic data. In *Proc. 37th Int. Conf. on Neural Inf. Process. Syst.*. Curran Associates Inc., Article 2208, 27 pages.
- [68] Tsu-Jui Fu, William Yang Wang, Daniel J. McDuff, and Yale Song. 2021. DOC2PPT: Automatic Presentation Slides Generation from Scientific Documents. In *AAAI Conf. on Artif. Intell.*.
- [69] Martin Funkquist, Iliia Kuznetsov, Yufang Hou, and Iryna Gurevych. 2023. CiteBench: A Benchmark for Scientific Citation Text Generation. In *Proc. 2023 Conf. on Empir. Methods Nat. Lang. Process.*. ACL, 7337–7353.
- [70] Catherine A Gao, Frederick M Howard, Nikolay S Markov, Emma C Dyer, Siddhi Ramesh, Yuan Luo, and Alexander T Pearson. 2023. Comparing scientific abstracts generated by ChatGPT to real abstracts with detectors and blinded human reviewers. *NPJ Digit. Med.* 6, 1 (2023), 75.
- [71] Andres Garcia-Silva, Cristian Berrio, and Jose Manuel Gomez-Perez. 2024. SPACE-IDEAS: A Dataset for Salient Information Detection in Space Innovation. In *Proc. 2024 Jt. Int. Conf. on Comput. Linguist. Lang. Resour. Eval. (LREC-COLING 2024)*. ELRA and ICCL, 15087–15092.
- [72] Eugene Garfield. 1955. Citation Indexes for Science. *Science* 122, 3159 (1955), 108–111.
- [73] Tirthankar Ghosal, Sandeep Kumar, Prabhat Kumar Bharti, and Asif Ekbal. 2022. Peer review analyze: A novel benchmark resource for computational analysis of peer reviews. *PLOS One* 17, 1 (01 2022), 1–29.
- [74] Tirthankar Ghosal, Kamal Kaushik Varanasi, and Valia Kordoni. 2022. HedgePeer: A dataset for uncertainty detection in peer reviews. In *Proc. 22nd ACM/IEEE Jt. Conf. on Digit. Libr.*. 1–5.
- [75] Tirthankar Ghosal, Rajeev Verma, Asif Ekbal, and Pushpak Bhattacharyya. 2019. DeepSentiPeer: Harnessing Sentiment in Review Texts to Recommend Peer Review Decisions. In *Proc. 57th Annu. Meet. Assoc. for Comput. Linguist.*. ACL, 1120–1130.
- [76] Hamed Babaei Giglou, Jennifer D’Souza, and Sören Auer. 2024. LLMs4Synthesis: Leveraging Large Language Models for Scientific Synthesis. In *2024 ACM/IEEE Jt. Conf. on Digit. Libr. (JCDL)*. 1–12.
- [77] Max Glockner, Yufang Hou, Preslav Nakov, and Iryna Gurevych. 2024. Grounding Fallacies Misrepresenting Scientific Publications in Evidence. [arXiv:2408.12812](https://arxiv.org/abs/2408.12812)
- [78] Max Glockner, Yufang Hou, Preslav Nakov, and Iryna Gurevych. 2024. Missci: Reconstructing Fallacies in Misrepresented Science. In *Proc. ACL*. ACL, 4372–4405.
- [79] Tomas Goldsack, Zhihao Zhang, Chenghua Lin, and Carolina Scarton. 2022. Making Science Simple: Corpora for the Lay Summarisation of Scientific Literature. In *Proc. 2022 Conf. on Empir. Methods Nat. Lang. Process.*. ACL, 10589–10604.
- [80] Andrew Gray. 2024. ChatGPT "contamination": estimating the prevalence of LLMs in the scholarly literature. [arXiv:2403.16887](https://arxiv.org/abs/2403.16887) [cs.DL]
- [81] Yanzhu Guo, Guokan Shang, Virgile Rennard, Michalis Vazirgiannis, and Chloé Clavel. 2023. Automatic Analysis of Substantiation in Scientific Peer Reviews. In *Find. Assoc. for Comput. Linguist. EMNLP 2023*. ACL, 10198–10216.
- [82] Thilo Hagendorff. 2024. Mapping the ethics of generative ai: A comprehensive scoping review. *arXiv preprint arXiv:2402.08323* (2024).
- [83] James Hartley. 2008. *Academic Writing and Publishing: A Practical Handbook*. Routledge.
- [84] Soheil HassaniPour, Sandeep Nayak, Ali Bozorgi, Mohammad-Hossein Keivanlou, Tirth Dave, Abdulhadi Alotaibi, Farahnaz Joukar, Parinaz Mellatdoust, Arash Bakhshi, Dona Kuriyakose, et al. 2024. The ability of ChatGPT in paraphrasing texts and reducing plagiarism: a descriptive analysis. *JMIR Med. Educ.* 10, 1 (2024), e53308.
- [85] Janna Hastings. 2023. *AI for Scientific Discovery*. CRC Press.
- [86] Xin He, Kaiyong Zhao, and Xiaowen Chu. 2021. AutoML: A survey of the state-of-the-art. *Knowledge-based systems* 212 (2021), 106622.
- [87] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. 2021. CLIPScore: A Reference-free Evaluation Metric for Image Captioning. In *Proc. 2021 Conf. on Empir. Methods Nat. Lang. Process.*. ACL, 7514–7528.
- [88] Tony Hey, Stewart Tansley, , and Kristin Tolle. 2009. Jim Gray on eScience: a Transformed Scientific Method. In *The Fourth Paradigm: Data-Intensive Scientific Discovery*, Tony Hey, Stewart Tansley, , and Kristin Tolle (Eds.). Microsoft Research.
- [89] Cong Duy Vu Hoang and Min-Yen Kan. 2010. Towards Automated Related Work Summarization. In *Coling 2010: Posters*, Chu-Ren Huang and Dan Jurafsky (Eds.). Coling 2010 Organizing Committee, Beijing, China, 427–435.
- [90] Sally Hopewell, Mike Clarke, David Moher, Elizabeth Wager, Philippa Middleton, Douglas G Altman, Kenneth F Schulz, and Consort Group. 2008. CONSORT for reporting randomized controlled trials in journal and conference abstracts: explanation and elaboration. *PLoS medicine* 5, 1 (2008), e20.
- [91] Yufang Hou, Charles Jochim, Martin Gleize, Francesca Bonin, and Debasis Ganguly. 2019. Identification of Tasks, Datasets, Evaluation Metrics, and Numeric Scores for Scientific Leaderboards Construction, Anna Korhonen, David Traum, and Lluís Màrquez (Eds.). ACL, Florence, Italy, 5203–5213.
- [92] Xiang Hu, Hongyu Fu, Jing Wang, Yifeng Wang, Zhikun Li, Renjun Xu, Yu Lu, Yaochu Jin, Lili Pan, and Zhenzhong Lan. 2024. Nova: An Iterative Planning and Search Approach to Enhance Novelty and Diversity of LLM Generated Ideas. *arXiv preprint arXiv:2410.14255* (2024).
- [93] Yue Hu and Xiaojun Wan. 2013. PPSGen: Learning to Generate Presentation Slides for Academic Papers. In *Int. Jt. Conf. on Artif. Intell.*.
- [94] Yue Hu and Xiaojun Wan. 2014. Automatic generation of related work sections in scientific papers: an optimization approach. In *Proc. 2014 Conf. on Empir. Methods Nat. Lang. Process. (EMNLP)*. 1624–1633.

- [95] Xinyu Hua, Mitko Nikolov, Nikhil Badugu, and Lu Wang. 2019. Argument Mining for Understanding Peer Reviews. In *Proc. 2019 Conf. North Am. Chapter Assoc. for Comput. Linguist. Hum. Lang. Technol. Vol. 1 (Long Short Pap.)*. ACL, 2131–2137.
- [96] Jie Huang and Kevin Chen-Chuan Chang. 2023. Citation: A key to building responsible and accountable large language models. *arXiv preprint arXiv:2307.02185* (2023).
- [97] Jingshan Huang and Ming Tan. 2023. The role of ChatGPT in scientific communication: writing better scientific review articles. *Am. journal cancer research* 13, 4 (2023), 1148.
- [98] Qian Huang, Jian Vora, Percy Liang, and Jure Leskovec. 2024. MAgentBench: Evaluating Language Agents on Machine Learning Experimentation. In *Forty-first Int. Conf. on Mach. Learn.*.
- [99] Taesoon Hwang, Nishant Aggarwal, Pir Zarak Khan, Thomas Roberts, Amir Mahmood, Madlen M Griffiths, Nick Parsons, and Saboor Khan. 2024. Can ChatGPT assist authors with abstract writing in medical journals? Evaluating the quality of scientific abstracts generated by ChatGPT and original abstracts. *Plos one* 19, 2 (2024), e0297701.
- [100] Joseph James, Chenghao Xiao, Yucheng Li, and Chenghua Lin. 2024. On the Rigour of Scientific Writing: Criteria, Analysis, and Insights. In *Find. Assoc. for Comput. Linguist. EMNLP 2024*. ACL, 6523–6538.
- [101] Peiwen Jiang, Xinbo Lin, Zibo Zhao, Ruhui Ma, Yvonne Jie Chen, and Jinhua Cheng. 2024. TKG: Redefinition and A New Way of Text-to-Table Tasks Based on Real World Demands and Knowledge Graphs Augmented LLMs. In *Proc. 2024 Conf. on Empir. Methods Nat. Lang. Process.*. ACL, 16112–16126.
- [102] Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. 2024. SWE-bench: Can Language Models Resolve Real-world Github Issues?. In *The Twelfth Int. Conf. on Learn. Represent.*.
- [103] Samira Ebrahimi Kahou, Adam Atkinson, Vincent Michalski, Ákos Kádár, Adam Trischler, and Yoshua Bengio. 2018. FigureQA: An Annotated Figure Dataset for Visual Reasoning.
- [104] Dongyeop Kang, Waleed Ammar, Bhavana Dalvi, Madeleine van Zuylen, Sebastian Kohlmeier, Eduard Hovy, and Roy Schwartz. 2018. A Dataset of Peer Reviews (PeerRead): Collection, Insights and NLP Applications. In *Proc. 2018 Conf. North Am. Chapter Assoc. for Comput. Linguist. Hum. Lang. Technol. Vol. 1 (Long Pap.)*. ACL, 1647–1661.
- [105] SeongKu Kang, Yunyi Zhang, Pengcheng Jiang, Dongha Lee, Jiawei Han, and Hwanjo Yu. 2024. Taxonomy-guided Semantic Indexing for Academic Paper Search. In *Proc. 2024 Conf. on Empir. Methods Nat. Lang. Process.*. ACL, 7169–7184.
- [106] Wei-Yu Kao and An-Zi Yen. 2024. How We Refute Claims: Automatic Fact-Checking through Flaw Identification and Explanation. In *Companion Proc. ACM on Web Conf. 2024*. 758–761.
- [107] Marcin Kardas, Piotr Czapla, Pontus Stenetorp, Sebastian Ruder, Sebastian Riedel, Ross Taylor, and Robert Stojnic. 2020. AxCell: Automatic Extraction of Results from Machine Learning Papers, Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu (Eds.). ACL, Online, 8580–8594.
- [108] Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali Farhadi. 2016. A Diagram is Worth a Dozen Images. In *Comput. Vis. – ECCV 2016*. Springer International Publishing, 235–251.
- [109] Neha Nayak Kennard, Tim O’Gorman, Rajarshi Das, Akshay Sharma, Chhandak Bagchi, Matthew Clinton, Pranay Kumar Yelugam, Hamed Zamani, and Andrew McCallum. 2022. DISAPERE: A Dataset for Discourse Structure in Peer Review Discussions. In *Proc. NAACL*. ACL, 1234–1249.
- [110] Seong-Gon Kim. 2023. Using ChatGPT for language editing in scientific articles. *Maxillofac. plastic reconstructive surgery* 45, 1 (2023), 13.
- [111] William R. King and Jun He. 2005. Understanding the Role and Methods of Meta-Analysis in IS Research. *Commun. Assoc. for Inf. Syst.* 16 (2005), 665–686.
- [112] Roger E. Kirk. 2009. Experimental Design. In *The SAGE Handbook of Quantitative Methods in Psychology*, Roger E. Millsap and Alberto Maydeu-Olivares (Eds.). 23–45.
- [113] Petr Knuth, Drahomira Herrmannova, Matteo Cancellieri, Lucas Anastasiou, Nancy Pontika, Samuel Pearce, Bikash Gyawali, and David Pride. 2023. CORE: a Global aggregation Service for Open access Papers. *Sci. Data* 10, 1 (2023), 366.
- [114] Dmitry Kobak, Rita González-Márquez, Emőke-Ágnes Horvát, and Jan Lause. 2024. Delving into ChatGPT usage in academic writing through excess vocabulary. *arXiv preprint arXiv:2406.07016* (2024).
- [115] Jing Yu Koh, Stephen McAleer, Daniel Fried, and Ruslan Salakhutdinov. 2024. Tree Search for Language Model Agents. *arXiv preprint arXiv:2407.01476* (2024).
- [116] Buse Sibel Korkmaz and Antonio Del Rio Chanona. 2024. Integrating Table Representations into Large Language Models for Improved Scholarly Document Comprehension. In *Proc. Fourth Workshop on Sch. Document Process. (SDP 2024)*. ACL, 293–306.
- [117] Kayvan Kousha and Mike Thelwall. 2024. Artificial intelligence to support publishing and peer review: A summary and review. *Learn. Publ.* 37, 1 (2024), 4–12.
- [118] Asheesh Kumar, Tirthankar Ghosal, and Asif Ekbal. 2021. A Deep Neural Architecture for Decision-Aware Meta-Review Generation. In *2021 ACM/IEEE Jt. Conf. on Digit. Libr. (JCDL)*. 222–225.
- [119] Nino Künzli, Anke Berger, Katarzyna Czabanowska, Raquel Lucas, Andrea Madarasova Geckova, Sarah Mantwill, and Olaf von Dem Knesebeck. 2022. «I Do Not Have Time» - Is This the End of Peer Review in Public Health Sciences? *Public Health Rev.* 43 (2022).
- [120] Iliia Kuznetsov, Osama Mohammed Afzal, Koen Dercksen, Nils Dycke, Alexander Goldberg, Tom Hope, Dirk Hovy, Jonathan K Kummerfeld, Anne Lauscher, Kevin Leyton-Brown, et al. 2024. What Can Natural Language Processing Do for Peer Review? *arXiv preprint arXiv:2405.06563* (2024).
- [121] Iliia Kuznetsov, Jan Buchmann, Max Eichler, and Iryna Gurevych. 2022. Revise and Resubmit: An Intertextual Model of Text-based Collaboration in Peer Review. *Comput. Linguist.* 48, 4 (12 2022), 949–986.

- [122] Pat Langley. 2000. The Computational Support of Scientific Discovery. *Int. J. Human-Computer Stud.* 53, 3 (2000), 393–410.
- [123] Po-Shen Lee, Jevin D. West, and Bill Howe. 2016. Viziometrics: Analyzing Visual Information in the Scientific Literature. *IEEE Trans. on Big Data* 4 (2016), 117–129.
- [124] Christoph Leiter, Jonas Belouadi, Yanran Chen, Ran Zhang, Daniil Larionov, Aida Kostikova, and Steffen Eger. 2024. NLLG Quarterly arXiv Report 09/24: What are the most influential current AI Papers? *ArXiv abs/2412.12121* (2024).
- [125] Adrian Letchford, Helen Susannah Moat, and Tobias Preis. 2015. The Advantage of Short Paper Titles. *Royal Soc. Open Sci.* 2, 8 (2015).
- [126] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive NLP tasks. In *Proc. NeurIPS*, Vol. 33. 9459–9474.
- [127] Da-Wei Li, Danqing Huang, Tingting Ma, and Chin-Yew Lin. 2021. Towards Topic-Aware Slide Generation For Academic Papers With Unsupervised Mutual Learning. In *AAAI Conf. on Artif. Intell.*
- [128] Junyi Li, Xiaoxue Cheng, Wayne Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. 2023. Halueval: A large-scale hallucination evaluation benchmark for large language models. *arXiv preprint arXiv:2305.11747* (2023).
- [129] Jiyi Li, Ayaka Sato, Kazuya Shimura, and Fumiyo Fukumoto. 2020. Multi-task Peer-Review Score Prediction. In *Proc. First Workshop on Sch. Document Process.* ACL, 121–126.
- [130] Lei Li, Yuqi Wang, Runxin Xu, Peiyi Wang, Xiachong Feng, Lingpeng Kong, and Qi Liu. 2024. Multimodal ArXiv: A Dataset for Improving Scientific Comprehension of Large Vision-Language Models. In *Proc. 62nd Annu. Meet. Assoc. for Comput. Linguist.*, Vol. 1. ACL, 14369–14387.
- [131] Long Li, Weiwen Xu, Jiayan Guo, Ruochen Zhao, Xingxuan Li, Yuqian Yuan, Boqiang Zhang, Yuming Jiang, Yifei Xin, Ronghao Dang, et al. 2024. Chain of Ideas: Revolutionizing Research Via Novel Idea Development with LLM Agents. *arXiv preprint arXiv:2410.13185* (2024).
- [132] Miao Li, Eduard Hovy, and Jey Lau. 2023. Summarizing Multiple Documents with Conversational Structure for Meta-Review Generation. In *Find. Assoc. for Comput. Linguist. EMNLP 2023*. ACL, 7089–7112.
- [133] Shujie Li, Liang Li, Ruiying Geng, Min Yang, Binhua Li, Guanghu Yuan, Wanwei He, Shao Yuan, Can Ma, Fei Huang, and Yongbin Li. 2024. Unifying Structured Data as Graph for Data-to-Text Pre-Training. *Trans. Assoc. for Comput. Linguist.* 12 (03 2024), 210–228. https://doi.org/10.1162/tacl_a_00641 arXiv:https://direct.mit.edu/tacl/article-pdf/doi/10.1162/tacl_a_00641/2346090/tacl_a_00641.pdf
- [134] Weitao Li, Junkai Li, Weizhi Ma, and Yang Liu. 2024. Citation-Enhanced Generation for LLM-based Chatbot. *arXiv preprint arXiv:2402.16063* (2024).
- [135] Xiangci Li, Yi-Hui Lee, and Jessica Ouyang. 2024. Cited Text Spans for Scientific Citation Text Generation. In *Proc. Fourth Workshop on Sch. Document Process. (SDP 2024)*. 90–104.
- [136] Xiangci Li, Biswadip Mandal, and Jessica Ouyang. 2022. CORWA: A Citation-Oriented Related Work Annotation Dataset. In *Proc. 2022 Conf. North Am. Chapter Assoc. for Comput. Linguist. Hum. Lang. Technol.*, Marine Carpuat, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz (Eds.). ACL, Seattle, United States, 5426–5440.
- [137] Xiangci Li and Jessica Ouyang. 2024. Related Work and Citation Text Generation: A Survey. In *Proc. 2024 Conf. on Empir. Methods Nat. Lang. Process.* ACL, 13846–13864.
- [138] Yutong Li, Lu Chen, Aiwei Liu, Kai Yu, and Lijie Wen. 2025. ChatCite: LLM Agent with Human Workflow Guidance for Comparative Literature Summary. In *Proc. 31st Int. Conf. on Comput. Linguist.* ACL, 3613–3630.
- [139] Zehan Li, Xin Zhang, Yanzhao Zhang, Dingkun Long, Pengjun Xie, and Meishan Zhang. 2023. Towards general text embeddings with multi-stage contrastive learning. *arXiv preprint arXiv:2308.03281* (2023).
- [140] Weixin Liang, Yuhui Zhang, Hancheng Cao, Binglu Wang, Daisy Ding, Xinyu Yang, Kailas Vodrahalli, Siyu He, Daniel Smith, Yian Yin, et al. 2023. Can large language models provide useful feedback on research papers? A large-scale empirical analysis. *arXiv preprint arXiv:2310.01783* (2023).
- [141] Weixin Liang, Yaohui Zhang, Zhengxuan Wu, Haley Lepp, Wenlong Ji, Xuandong Zhao, Hancheng Cao, Sheng Liu, Siyu He, Zhi Huang, Diyi Yang, Christopher Potts, Christopher D Manning, and James Y. Zou. 2024. Mapping the Increasing Use of LLMs in Scientific Papers. In *Proc. COLM*.
- [142] Jialiang Lin, Jiaxin Song, Zhangping Zhou, Yidong Chen, and Xiaodong Shi. 2023. Automated scholarly paper review: Concepts, technologies, and challenges. *Inf. Fusion* 98 (2023), 101830.
- [143] Jialiang Lin, Jiaxin Song, Zhangping Zhou, Yidong Chen, and Xiaodong Shi. 2023. MOPRD: A multidisciplinary open peer review dataset. *Neural Comput. Appl.* 35, 34 (Sept. 2023), 24191–24206.
- [144] Stephanie Lin, Jacob Hilton, and Owain Evans. 2021. Truthfulqa: Measuring how models mimic human falsehoods. *arXiv preprint arXiv:2109.07958* (2021).
- [145] Can Liu, Yun Han, Ruikang Jiang, and Xiaoru Yuan. 2021. Advisor: Automatic visualization answer for natural-language question on tabular data. In *2021 IEEE 14th Pac. Vis. Symp. (PacificVis)*. IEEE, 11–20.
- [146] Cencen Liu, Yi Xu, Wen Yin, and Dezhong Zheng. 2023. Structure-aware Table-to-Text Generation with Prefix-tuning. In *Proc. 2023 4th Int. Conf. on Control. Robotics Intell. Syst.* Association for Computing Machinery, 135–140.
- [147] Haokun Liu, Yangqiaoyu Zhou, Mingxuan Li, Chenfei Yuan, and Chenhao Tan. 2024. Literature meets data: A synergistic approach to hypothesis generation. *arXiv preprint arXiv:2410.17309* (2024).
- [148] Quanliang Liu, Maciej P Polak, So Yeon Kim, MD Shuvo, Hrishikesh Shridhar Deodhar, Jeongsoo Han, Dane Morgan, and Hyunseok Oh. 2024. Beyond designer’s knowledge: Generating materials design hypotheses via large language models. *arXiv preprint arXiv:2409.06756* (2024).
- [149] Ryan Liu and Nihar B. Shah. 2023. ReviewerGPT? An Exploratory Study on Using Large Language Models for Paper Reviewing. *arXiv:2306.00622* (2023).

- [150] Sizhe Liu, Yizhou Lu, Siyu Chen, Xiyang Hu, Jieyu Zhao, Tianfan Fu, and Yue Zhao. 2024. DrugAgent: Automating AI-aided Drug Discovery Programming through LLM Multi-Agent Collaboration. *arXiv preprint arXiv:2411.15692* (2024).
- [151] Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. 2024. The AI Scientist: Towards Fully Automated Open-Ended Scientific Discovery. *arXiv preprint arXiv:2408.06292* (2024).
- [152] Yuyu Luo, Jiawei Tang, and Guoliang Li. 2021. nvBench: A Large-Scale Synthesized Dataset for Cross-Domain Natural Language to Visualization Task. *ArXiv abs/2112.12926* (2021).
- [153] Ziming Luo, Zonglin Yang, Zexin Xu, Wei Yang, and Xinya Du. 2025. LLM4SR: A Survey on Large Language Models for Scientific Research. *arXiv:2501.04306* [cs.CL]
- [154] Kangyong Ma. 2025. AI agents in chemical research: GVIM – an intelligent research assistant system. *Digit. Discov.* (jan 2025).
- [155] Calum Macdonald, Davies Adeloye, Aziz Sheikh, and Igor Rudan. 2023. Can ChatGPT draft a research article? An example of population-level vaccine effectiveness analysis. *J. global health* 13 (2023).
- [156] Jock Mackinlay. 1986. Automating the design of graphical presentations of relational information. *Acm Trans. On Graph. (Tog)* 5, 2 (1986), 110–141.
- [157] Paula Maddigan and Teo Susnjak. 2023. Chat2VIS: Generating Data Visualizations via Natural Language Using ChatGPT, Codex and GPT-3 Large Language Models. *IEEE Access* 11 (2023), 45181–45193.
- [158] Himanshu Maheshwari, Sambaran Bandyopadhyay, Aparna Garimella, and Anandhavelu Natarajan. 2024. Presentations are not always linear! GNN meets LLM for Text Document-to-Presentation Transformation with Attribution. In *Find. Assoc. for Comput. Linguist. EMNLP 2024. ACL*, 15948–15962.
- [159] Ahmed Masry, Xuan Long Do, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. 2022. ChartQA: A Benchmark for Question Answering about Charts with Visual and Logical Reasoning. In *Find. Assoc. for Comput. Linguist. ACL 2022. ACL*, 2263–2279.
- [160] Prakhar Mishra, Chaitali Diwan, Srinath Srinivasa, and Gopalakrishnan Srinivasaraghavan. 2021. Automatic title generation for text with pre-trained transformer language model. In *2021 IEEE 15th Int. Conf. on Semantic Comput. (ICSC)*. IEEE, 17–24.
- [161] Ishani Mondal, Zongxia Li, Yufang Hou, Anandhavelu Natarajan, Aparna Garimella, and Jordan Lee Boyd-Graber. 2024. SciDoc2Diagrammer-MAF: Towards Generation of Scientific Diagrams from Documents guided by Multi-Aspect Feedback Refinement. In *Find. Assoc. for Comput. Linguist. EMNLP 2024. ACL*, 13342–13375.
- [162] Ishani Mondal, Shwetha S, Anandhavelu Natarajan, Aparna Garimella, Sambaran Bandyopadhyay, and Jordan Boyd-Graber. 2024. Presentations by the Humans and For the Humans: Harnessing LLMs for Generating Persona-Aware Slides from Documents. In *Proc. EACL. ACL*, 2664–2684.
- [163] Nafise Moosavi, Andreas Rücklé, Dan Roth, and Iryna Gurevych. 2021. SciGen: a Dataset for Reasoning-Aware Text Generation from Scientific Tables. In *Proc. Neural Inf. Process. Syst. Track on Datasets Benchmarks*, Vol. 1.
- [164] Niklas Muennighoff, Nouamane Tazi, Loic Magne, and Nils Reimers. 2023. MTEB: Massive Text Embedding Benchmark. In *Proc. 17th Conf. Eur. Chapter Assoc. for Comput. Linguist.. ACL*, 2014–2037.
- [165] Arpit Narechania, Arjun Srinivasan, and John Stasko. 2020. NL4DV: A toolkit for generating analytic specifications for data visualization from natural language queries. *IEEE Trans. on Vis. Comput. Graph.* 27, 2 (2020), 369–379.
- [166] Deepak Nathani, Lovish Madaan, Nicholas Roberts, Nikolay Bashlykov, Ajay Menon, Vincent Moens, Amar Budhiraja, Despoina Magka, Vladislav Vorotilov, Gaurav Chaurasia, Dieuwke Hupkes, Ricardo Silveira Cabral, Tatiana Shavrina, Jakob Foerster, Yoram Bachrach, William Yang Wang, and Roberta Raileanu. 2025. MLGym: A New Framework and Benchmark for Advancing AI Research Agents. *arXiv:2502.14499* [cs.CL]
- [167] Benjamin Newman, Yoonjoo Lee, Aakanksha Naik, Pao Siangliulue, Raymond Fok, Juho Kim, Daniel S Weld, Joseph Chee Chang, and Kyle Lo. 2024. ArxivDIGESTables: Synthesizing Scientific Literature into Tables using Language Models. In *Proc. 2024 Conf. on Empir. Methods Nat. Lang. Process.. ACL*, 9612–9631.
- [168] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proc. 40th annual meeting Assoc. for Comput. Linguist.*, 311–318.
- [169] Ankur Parikh, Xuezhi Wang, Sebastian Gehrmann, Manaal Faruqui, Bhuwan Dhingra, Diyi Yang, and Dipanjan Das. 2020. ToTTo: A Controlled Table-To-Text Generation Dataset. In *Proc. 2020 Conf. on Empir. Methods Nat. Lang. Process. (EMNLP)*. ACL, 1173–1186.
- [170] Yang Jeong Park, Daniel Kaplan, Zhichu Ren, Chia-Wei Hsu, Changhao Li, Haowei Xu, Sipei Li, and Ju Li. 2024. Can ChatGPT be used to generate scientific hypotheses? *J. Materiomics* 10, 3 (2024), 578–584.
- [171] Guy Paré, Marie-Claude Trudel, Mirou Jaana, and Spyros Kitsiou. 2015. Synthesizing Information Systems Knowledge: a Typology of Literature Reviews. *Inf. & Manag.* 52 (2015), 183–199.
- [172] Dominic Petrak, Nafise Sadat Moosavi, and Iryna Gurevych. 2023. Arithmetic-Based Pretraining Improving Numeracy of Pretrained Language Models. In *Proc. 12th Jt. Conf. on Lexical Comput. Semant. (*SEM 2023)*. ACL, 477–493.
- [173] Nicky Phillips. 2017. Online software spots genetic errors in cancer papers. *Nature* 551, 7681 (2017).
- [174] Barbara Plank and Reinard van Dalen. 2019. CiteTracked: A longitudinal dataset of peer reviews and citations. In *Proc. 4th Jt. Workshop on Bibliometric-enhanced Inf. Retr. Nat. Lang. Process. for Digit. Libr. (BIRNDL 2019) co-located with 42nd Int. ACM SIGIR Conf. on Res. Dev. Inf. Retr. (SIGIR 2019)*. CEUR Workshop Proceedings, 116–122.
- [175] Aniket Pramanick, Yufang Hou, Saif M. Mohammad, and Iryna Gurevych. 2024. The Nature of NLP: Analyzing Contributions in NLP Papers. *arXiv:2409.19505* [cs.CL]
- [176] Aniket Pramanick, Yufang Hou, Saif M. Mohammad, and Iryna Gurevych. 2024. Transforming Scholarly Landscapes: Influence of Large Language Models on Academic Fields beyond Computer Science. *arXiv:2409.19508* [cs.CL]

- [177] Shraman Pramanick, Rama Chellappa, and Subhashini Venugopalan. 2024. SPIQA: A Dataset for Multimodal Question Answering on Scientific Papers. In *The Thirty-eight Conf. on Neural Inf. Process. Syst. Datasets Benchmarks Track*.
- [178] Kevin Pu, KJ Feng, Tovi Grossman, Tom Hope, Bhavana Dalvi Mishra, Matt Latzke, Jonathan Bragg, Joseph Chee Chang, and Pao Siangliulue. 2024. IdeaSynth: Iterative Research Idea Development Through Evolving and Composing Idea Facets with Literature-Grounded Feedback. *arXiv preprint arXiv:2410.04025* (2024).
- [179] Zhongji Pu, Chun-Lin Shi, Che Ok Jeon, Jingyuan Fu, Shuang-Jiang Liu, Canhui Lan, Yanlai Yao, Yong-Xin Liu, and Baolei Jia. 2024. ChatGPT and generative AI are revolutionizing the scientific community: A Janus-faced conundrum. *iMeta* 3, 2 (2024), e178.
- [180] Sukannya Purkayastha, Anne Lauscher, and Iryna Gurevych. 2023. Exploring Jiu-Jitsu Argumentation for Writing Peer Review Rebuttals. In *Proc. 2023 Conf. on Empir. Methods Nat. Lang. Process.*. ACL, 14479–14495.
- [181] Biqing Qi, Kaiyan Zhang, Haoxiang Li, Kai Tian, Sihang Zeng, Zhang-Ren Chen, and Bowen Zhou. 2023. Large language models are zero shot hypothesis proposers. *arXiv preprint arXiv:2311.05965* (2023).
- [182] Biqing Qi, Kaiyan Zhang, Kai Tian, Haoxiang Li, Zhang-Ren Chen, Sihang Zeng, Ermo Hua, Hu Jinfang, and Bowen Zhou. 2024. Large Language Models as Biomedical Hypothesis Generators: A Comprehensive Evaluation. *arXiv preprint arXiv:2407.08940* (2024).
- [183] Yuting Qiang, Yanwei Fu, Yanwen Guo, Zhi-Hua Zhou, and Leonid Sigal. 2016. Learning to Generate Posters of Scientific Papers. In *AAAI Conf. on Artif. Intell.*.
- [184] Marissa Radensky, Simra Shahid, Raymond Fok, Pao Siangliulue, Tom Hope, and Daniel S Weld. 2024. Scideator: Human-LLM Scientific Idea Generation Grounded in Research-Paper Facet Recombination. *arXiv preprint arXiv:2409.14634* (2024).
- [185] Raian Rahman, Rizvi Hasan, Abdullah Al Farhad, Md Tahmid Rahman Laskar, Md. Hamjajul Ashmafee, and Abu Raihan Mostofa Kamal. 2023. ChartSumm: A Comprehensive Benchmark for Automatic Chart Summarization of Long and Short Summaries. *ArXiv abs/2304.13620* (2023).
- [186] Pritika Ramu, Aparna Garimella, and Sambaran Bandyopadhyay. 2024. Is This a Bad Table? A Closer Look at the Evaluation of Table Generation from Text. In *Proc. 2024 Conf. on Empir. Methods Nat. Lang. Process.*. ACL, 22206–22216.
- [187] Chandan K Reddy and Parshin Shojaei. 2024. Towards Scientific Discovery with Generative AI: Progress, Opportunities, and Challenges. *arXiv:2412.11427* [cs.LG]
- [188] Isabelle Régner, Catherine Thinus-Blanc, Agnès Netter, Toni Schmader, and Pascal Huguet. 2019. Committees with implicit biases promote fewer women when they do not believe gender bias exists. *Nat. Hum. Behav.* 3, 11 (2019), 1171–1179.
- [189] Georg Rehm, Stefan Dietze, Sonja Schimmler, and Frank Krüger. 2024. *Natural Scientific Language Processing and Research Knowledge Graphs: First International Workshop, NSLP 2024, Hersonissos, Crete, Greece, May 27, 2024, Proceedings*. Number 14770. Springer.
- [190] Zachary Robertson. 2023. GPT-4 is Slightly Helpful for Peer-Review Assistance: A Pilot Study. *arXiv preprint arXiv:2307.05492* (2023).
- [191] Steven F Roth, John Kolojchick, Joe Mattis, and Jade Goldstein. 1994. Interactive graphic design using automatic presentation knowledge. In *Proc. SIGCHI conference on Hum. factors computing systems*. 112–117.
- [192] Steven F Roth and Joe Mattis. 1991. Automating the presentation of information. In *Proc. IEEE Conf. on Artif. Intell. Appl.*. 90–97.
- [193] Furkan Şahinuç, Ilia Kuznetsov, Yufang Hou, and Iryna Gurevych. 2024. Systematic Task Exploration with LLMs: A Study in Citation Text Generation. In *Proc. 62nd Annu. Meet. Assoc. for Comput. Linguist.*, Vol. 1. ACL, 4832–4855.
- [194] Furkan Şahinuç, Thy Thy Tran, Yulia Grishina, Yufang Hou, Bei Chen, and Iryna Gurevych. 2024. Efficient Performance Tracking: Leveraging Large Language Models for Automated Construction of Scientific Leaderboards. In *Proc. 2024 Conf. on Empir. Methods Nat. Lang. Process.*. ACL, 7963–7977.
- [195] Michele Salvagno, Fabio Silvio Taccone, and Alberto Giovanni Gerli. 2023. Can artificial intelligence help for scientific writing? *Crit. care* 27, 1 (2023), 75.
- [196] Adam Santoro, David Raposo, David G Barrett, Mateusz Malinowski, Razvan Pascanu, Peter Battaglia, and Timothy Lillicrap. 2017. A simple neural network module for relational reasoning. In *Adv. Neural Inf. Process. Syst.*, Vol. 30. Curran Associates, Inc.
- [197] Shubhra Kanti Karmaker Santu, Sanjeev Kumar Sinha, Naman Bansal, Alex Knipper, Souvika Sarkar, John Salvador, Yash Mahajan, Sri Guttikonda, Mousumi Akter, Matthew Freestone, and Matthew C. Williams Jr. 2024. Prompting LLMs to Compose Meta-Review Drafts from Peer-Review Narratives of Scholarly Manuscripts. *arXiv preprint arXiv:2402.15589* (2024). *arXiv:2402.15589* [cs.CL]
- [198] Laurie A. Schintler, Connie L. McNeely, and James Witte. 2023. A Critical Examination of the Ethics of AI-Mediated Peer Review. *arXiv preprint arXiv:2309.12356* (2023).
- [199] Daniel Schlagwein and Leslie Willcocks. 2023. ‘ChatGPT et al.’: The ethics of using (generative) artificial intelligence in research and science. *J. Inf. Technol.* 38, 3 (2023), 232–238.
- [200] Michael Schlichtkrull, Nedjma Ousidhoum, and Andreas Vlachos. 2023. The intended uses of automated fact-checking artefacts: Why, how and who. *arXiv preprint arXiv:2304.14238* (2023).
- [201] Dominik Schmidt, Zhengyao Jiang, and Yuxiang Wu. 2024. Introducing Weco AIDE. <https://www.weco.ai/blog/technical-report>.
- [202] SciScore. 2024. The best methods review tool for scientific research. <https://sciscore.com/>. Accessed: 12 June 2024.
- [203] Chenhui Shen, Liying Cheng, Ran Zhou, Lidong Bing, Yang You, and Luo Si. 2022. MReD: A Meta-Review Dataset for Structure-Controllable Text Generation. In *Find. Assoc. for Comput. Linguist. ACL 2022*. ACL, 2521–2535.
- [204] Chufan Shi, Cheng Yang, Yaxin Liu, Bo Shui, Junjie Wang, Mohan Jing, Linran Xu, Xinyu Zhu, Siheng Li, Yuxiang Zhang, Gongye Liu, Xi-aomei Nie, Deng Cai, and Yujiu Yang. 2024. ChartMimic: Evaluating LLM’s Cross-Modal Reasoning Capability via Chart-to-Code Generation. *arXiv:2406.09961* [cs.SE]

- [205] Haoxiang Shi, Jiaan Wang, Jiarong Xu, Cen Wang, and Tetsuya Sakai. 2024. CT-Eval: Benchmarking Chinese Text-to-Table Performance in Large Language Models. *arXiv preprint arXiv:2405.12174* (2024).
- [206] Chenglei Si, Diyi Yang, and Tatsunori Hashimoto. 2024. Can LLMs Generate Novel Research Ideas? A Large-Scale Human Study with 100+ NLP Researchers. *arXiv:2409.04109*
- [207] Shruti Singh, Mayank Singh, and Pawan Goyal. 2021. COMPARE: A Taxonomy and Dataset of Comparison Discussions in Peer Reviews. *arXiv:2108.04366* [cs.CL]
- [208] Wael Soliman and Mikko Siponen. 2022. What Do We Really Mean by Rigor in Information Systems Research?
- [209] Yuanfeng Song, Xuefang Zhao, Raymond Chi-Wing Wong, and Di Jiang. 2022. RGVisNet: A Hybrid Retrieval-Generation Neural Framework Towards Automatic Data Visualization Generation. In *Proc. 28th ACM SIGKDD Conf. on Knowl. Discov. Data Min.*. Association for Computing Machinery, 1646–1655.
- [210] Karen Sparck Jones. 1972. A statistical interpretation of term specificity and its application in retrieval. *J. documentation* 28, 1 (1972), 11–21.
- [211] M. Sravanthi, C. Ravindranath Chowdary, and P. Sreenivasa Kumar. 2009. SlidesGen: Automatic Generation of Presentation Slides for a Technical Paper Using Summarization. In *The Fla. AI Res. Soc.*
- [212] Lukas Stappen, Georgios Rizos, Madina Hasan, Thomas Hain, and Björn W Schuller. 2020. Uncertainty-aware machine support for paper reviewing on the interspeech 2019 submission corpus. (2020).
- [213] Moritz Staudinger, Wojciech Kusa, Florina Piroi, and Allan Hanbury. 2024. An Analysis of Tasks and Datasets in Peer Reviewing. In *Proc. Fourth Workshop on Sch. Document Process. (SDP 2024)*. ACL, 257–268.
- [214] Dana Strauss, Sophia Gran-Ruaz, Muna Osman, Monnica T Williams, and Sonya C Faber. 2023. Racism and censorship in the editorial and peer review process. *Front. Psychol.* 14 (2023).
- [215] Haoyang Su, Renqi Chen, Shixiang Tang, Xinzhe Zheng, Jingzhe Li, Zhenfei Yin, Wanli Ouyang, and Nanqing Dong. 2024. Two Heads Are Better Than One: A Multi-Agent System Has the Potential to Improve Scientific Idea Generation. *arXiv preprint arXiv:2410.09403* (2024).
- [216] Lya Hulliyyatus Suadaa, Hidetaka Kamigaito, Kotaro Funakoshi, Manabu Okumura, and Hiroya Takamura. 2021. Towards Table-to-Text Generation with Numerical Reasoning. In *Proc. 59th Annu. Meet. Assoc. for Comput. Linguist. 11th Int. Jt. Conf. on Nat. Lang. Process.*, Vol. 1. ACL, 1451–1465.
- [217] Edward Sun, Yufang Hou, Dakuo Wang, Yunfeng Zhang, and Nancy X. R. Wang. 2021. D2S: Document-to-Slide Generation Via Query-Based Text Summarization. In *Proc. NAACL. ACL*, 1405–1418.
- [218] Kexuan Sun, Zhiqiang Qiu, Abel Salinas, Yuzhong Huang, Dong-Ho Lee, Daniel Benjamin, Fred Morstatter, Xiang Ren, Kristina Lerman, and Jay Pujara. 2022. Assessing scientific research papers with knowledge graphs. In *Proc. 45th Int. ACM SIGIR Conf. on Res. Dev. Inf. Retr.*. 2467–2472.
- [219] Lichao Sun, Yue Huang, Haoran Wang, Siyuan Wu, Qihui Zhang, Chujie Gao, Yixin Huang, Wenhan Lyu, Yixuan Zhang, Xiner Li, et al. 2024. Trustllm: Trustworthiness in large language models. *arXiv:2401.05561* (2024).
- [220] Yiwen Sun, Jason Leigh, Andrew Johnson, and Barbara Di Eugenio. 2014. Articulate: Creating meaningful visualizations from natural language. In *Innovative Approaches of Data Visualization and Visual Analytics*, Mao Lin Huang and Weidong Huang (Eds.). IGI Global, 218–235.
- [221] Anirudh Sundar, Christopher Richardson, and Larry Heck. 2024. gTBLS: Generating Tables from Text by Conditional Question Answering. *arXiv preprint arXiv:2403.14457* (2024).
- [222] Jaroslaw Szumega, Lamine Bougueroua, Blerina Gkotse, Pierre Jouvelot, and Federico Ravotti. 2023. The Open Review-Based (ORB) dataset: Towards Automatic Assessment of Scientific Papers and Experiment Proposals in High-Energy Physics. *arXiv:2312.04576* [cs.DL]
- [223] Neset Tan, Trung Nguyen, Josh Bensemann, Alex Peng, Qiming Bao, Yang Chen, Mark Gahegan, and Michael Witbrock. 2023. Multi2Claim: Generating Scientific Claims from Multi-Choice Questions for Scientific Fact-Checking. In *Proc. 17th Conf. Eur. Chapter Assoc. for Comput. Linguist.*. ACL, 2652–2664.
- [224] Jiawei Tang, Yuyu Luo, Mourad Ouzzani, Guoliang Li, and Hongyang Chen. 2022. Sevi: Speech-to-Visualization through Neural Machine Translation. In *Proc. 2022 Int. Conf. on Manag. Data*. Association for Computing Machinery, 2353–2356.
- [225] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288* (2023).
- [226] Patara Trirat, Wonyong Jeong, and Sung Ju Hwang. 2024. AutoML-Agent: A Multi-Agent LLM Framework for Full-Pipeline AutoML. *arXiv preprint arXiv:2410.02958* (2024).
- [227] Yun-Da Tsai, Yu-Che Tsai, Bo-Wei Huang, Chun-Pai Yang, and Shou-De Lin. 2023. Automl-gpt: Large language model for automl. *arXiv preprint arXiv:2309.01125* (2023).
- [228] Juraj Vladika and Florian Matthes. 2023. Scientific fact-checking: A survey of resources and approaches. *arXiv preprint arXiv:2305.16859* (2023).
- [229] Henrik Voigt, Kai Lawonn, and Sina Zarriß. 2024. Plots Made Quickly: An Efficient Approach for Generating Visualizations from Natural Language Queries. In *Proc. 2024 Jt. Int. Conf. on Comput. Linguist. Lang. Resour. Eval. (LREC-COLING 2024)*. 12787–12793.
- [230] Dario von Wedel, Rico A. Schmitt, Moritz Thiele, Raphael Leuner, Denys Shay, Simone Redaelli, and Maximilian S. Schaefer. 2024. Affiliation Bias in Peer Review of Abstracts by a Large Language Model. *JAMA* 331, 3 (01 2024), 252–253.
- [231] Tu Vu, Mohit Iyyer, Xuezhi Wang, Noah Constant, Jerry Wei, Jason Wei, Chris Tar, Yun-Hsuan Sung, Denny Zhou, Quoc Le, et al. 2023. Freshllms: Refreshing large language models with search engine augmentation. *arXiv preprint arXiv:2310.03214* (2023).
- [232] David Wadden, Kyle Lo, Bailey Kuehl, Arman Cohan, Iz Beltagy, Lucy Lu Wang, and Hannaneh Hajishirzi. 2022. SciFact-open: Towards open-domain scientific claim verification. *arXiv preprint arXiv:2210.13777* (2022).

- [233] William H Walters and Esther Isabelle Wilder. 2023. Fabrication and errors in the bibliographic citations generated by ChatGPT. *Sci. Reports* 13, 1 (2023), 14045.
- [234] Fei Wang, Zhewei Xu, Pedro Szekely, and Muhao Chen. 2022. Robust (Controlled) Table-to-Text Generation with Structure-Aware Equivariance Learning. In *Proc. 2022 Conf. North Am. Chapter Assoc. for Comput. Linguist. Hum. Lang. Technol.*. ACL, 5037–5048.
- [235] Hanchen Wang, Tianfan Fu, Yuanqi Du, Wenhao Gao, Kexin Huang, Ziming Liu, Payal Chandak, Shengchao Liu, Peter Van Katwyk, Andreea Deac, et al. 2023. Scientific discovery in the age of artificial intelligence. *Nature* 620, 7972 (2023), 47–60.
- [236] Hao Wang and Wu-Jun Li. 2014. Relational collaborative topic regression for recommender systems. *IEEE Trans. on Knowl. Data Eng.* 27, 5 (2014), 1343–1355.
- [237] Pancheng Wang, Shasha Li, Haifang Zhou, Jintao Tang, and Ting Wang. 2019. ToC-RWG: Explore the Combination of Topic Model and Citation Information for Automatic Related Work Generation. *IEEE Access* 8 (2019), 13043–13055.
- [238] Qingyun Wang, Lifu Huang, Zhiying Jiang, Kevin Knight, Heng Ji, Mohit Bansal, and Yi Luan. 2019. PaperRobot: Incremental Draft Generation of Scientific Ideas. In *Proc. 57th Annu. Meet. Assoc. for Comput. Linguist.*. ACL, 1980–1991.
- [239] Qingyun Wang, Qi Zeng, Lifu Huang, Kevin Knight, Heng Ji, and Nazneen Fatema Rajani. 2020. ReviewRobot: Explainable Paper Review Generation based on Knowledge Synthesis. In *Proc. 13th Int. Conf. on Nat. Lang. Gener.*. ACL, 384–397.
- [240] Sida Wang, Xiaojun Wan, and Shikang Du. 2017. Phrase-Based Presentation Slides Generation for Academic Papers. In *AAAI Conf. on Artif. Intell.*
- [241] Xingyao Wang, Boxuan Li, Yufan Song, Frank F Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan, Yueqi Song, Bowen Li, Jaskirat Singh, et al. 2024. OpenHands: An Open Platform for AI Software Developers as Generalist Agents. *arXiv preprint arXiv:2407.16741* (2024).
- [242] Yongzhen Wang, Xiaozhong Liu, and Zheng Gao. 2018. Neural Related Work Summarization with a Joint Context-driven Attention Mechanism. In *Proc. 2018 Conf. on Empir. Methods Nat. Lang. Process.*, Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun'ichi Tsujii (Eds.). ACL, Brussels, Belgium, 1776–1786.
- [243] 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. 2024. CharXiv: Charting Gaps in Realistic Chart Understanding in Multimodal LLMs. In *The Thirty-eight Conf. on Neural Inf. Process. Syst. Datasets Benchmarks Track*.
- [244] Ann C. Weller. 2001. *Editorial Peer Review: Its Strengths and Weaknesses*. American Society for Information Science and Technology.
- [245] Siwei Wu, Yizhi Li, Xingwei Qu, Rishi Ravikumar, Yucheng Li, Tyler Loakman Shanghaoran Quan Xiaoyong Wei, Riza Batista-Navarro, and Chenghua Lin. 2025. LongEval: A Comprehensive Analysis of Long-Text Generation Through a Plan-based Paradigm. *arXiv:2502.19103* [cs.CL]
- [246] Siwei Wu, Yizhi Li, Kang Zhu, Ge Zhang, Yiming Liang, Kaijing Ma, Chenghao Xiao, Haoran Zhang, Bohao Yang, Wenhua Chen, Wenhao Huang, Noura Al Moubayed, Jie Fu, and Chenghua Lin. 2024. SciMMIR: Benchmarking Scientific Multi-modal Information Retrieval. In *Find. Assoc. for Comput. Linguist. ACL 2024*. ACL, 12560–12574.
- [247] Guangzhi Xiong, Eric Xie, Amir Hassan Shariatmadari, Sikun Guo, Stefan Bekiranov, and Aidong Zhang. 2024. Improving Scientific Hypothesis Generation with Knowledge Grounded Large Language Models. *arXiv preprint arXiv:2411.02382* (2024).
- [248] Peixin Xu, Yajuan Ding, and Wenqi Fan. 2024. ChartAdapter: Large Vision-Language Model for Chart Summarization. *arXiv:2412.20715* [cs.MM]
- [249] Ruoxi Xu, Yingfei Sun, Mengjie Ren, Shiguang Guo, Ruotong Pan, Hongyu Lin, Le Sun, and Xianpei Han. 2024. AI for social science and social science of AI: A survey. *Inf. Process. & Manag.* 61, 3 (2024), 103665.
- [250] Sheng Xu and Xiaojun Wan. 2022. PosterBot: A System for Generating Posters of Scientific Papers with Neural Models. In *AAAI Conf. on Artif. Intell.*
- [251] Hidekazu Yanagimoto, Iroha Kisaku, and Kiyota Hashimoto. 2024. Table-to-Text Using Pre-trained Large Language Model and LoRA. In *2024 16th IIAI Int. Congr. on Adv. Appl. Informatics (IIAI-AAI)*. 91–96.
- [252] Zonglin Yang, Xinya Du, Junxian Li, Jie Zheng, Soujanya Poria, and Erik Cambria. 2023. Large language models for automated open-domain scientific hypotheses discovery. *arXiv preprint arXiv:2309.02726* (2023).
- [253] Zonglin Yang, Wanhao Liu, Ben Gao, Tong Xie, Yuqiang Li, Wanli Ouyang, Soujanya Poria, Erik Cambria, and Dongzhan Zhou. 2024. MOOSE-Chem: Large Language Models for Rediscovering Unseen Chemistry Scientific Hypotheses. *arXiv preprint arXiv:2410.07076* (2024).
- [254] Michihiro Yasunaga, Jungo Kasai, Rui Zhang, Alexander R Fabbri, Irene Li, Dan Friedman, and Dragomir R Radev. 2019. Scisummnet: A large annotated corpus and content-impact models for scientific paper summarization with citation networks. In *Proc. AAAI conference on artificial intelligence*, Vol. 33. 7386–7393.
- [255] Xinyang Yi, Ji Yang, Lichan Hong, Derek Zhiyuan Cheng, Lukasz Heldt, Aditee Kumthekar, Zhe Zhao, Li Wei, and Ed Chi. 2019. Sampling-bias-corrected neural modeling for large corpus item recommendations. In *Proc. 13th ACM Conf. on Recomm. Syst.*. Association for Computing Machinery, 269–277.
- [256] Zhangye Yin, Qiushi Sun, Qipeng Guo, Jiawen Wu, Xipeng Qiu, and Xuanjing Huang. 2023. Do Large Language Models Know What They Don't Know? *arXiv preprint arXiv:2305.18153* (2023).
- [257] Larry D. Yore, Brian M. Hand, and Marilyn K. Florence. 2004. Scientists' Views of Science, Models of Writing, and Science Writing Practices. *J. Res. Sci. Teach.* 41, 4 (2004), 338–369.
- [258] Yantao Yu, Weipeng Wang, Zhoutian Feng, and Daiyue Xue. 2021. A dual augmented two-tower model for online large-scale recommendation. *DLP-KDD* (2021).
- [259] Weizhe Yuan, Pengfei Liu, and Graham Neubig. 2022. Can We Automate Scientific Reviewing? *J. Artif. Int. Res.* 75 (dec 2022), 42 pages.

- [260] Abhay Zala, Han Lin, Jaemin Cho, and Mohit Bansal. 2024. DiagrammerGPT: Generating Open-Domain, Open-Platform Diagrams via LLM Planning. In *COLM*.
- [261] Qi Zeng, Mankeerat Sidhu, Hou Pong Chan, Lu Wang, and Heng Ji. 2023. Meta-review Generation with Checklist-guided Iterative Introspection. *arXiv preprint arXiv:2305.14647* (2023).
- [262] Jiajie Zhang, Zhongni Hou, Xin Lv, Shulin Cao, Zhenyu Hou, Yilin Niu, Lei Hou, Yuxiao Dong, Ling Feng, and Juanzi Li. 2024. LongReward: Improving Long-context Large Language Models with AI Feedback. *arXiv:2410.21252* [cs.CL]
- [263] Leixin Zhang, Steffen Eger, Yinjie Cheng, Weihe Zhai, Jonas Belouadi, Christoph Leiter, Simone Paolo Ponzetto, Fahimeh Moafian, and Zhixue Zhao. 2024. SciImage: How Good Are Multimodal Large Language Models at Scientific Text-to-Image Generation? *arXiv:2412.02368* [cs.AI]
- [264] Lei Zhang, Yuge Zhang, Kan Ren, Dongsheng Li, and Yuqing Yang. 2023. Mlcopilot: Unleashing the power of large language models in solving machine learning tasks. *arXiv preprint arXiv:2304.14979* (2023).
- [265] Xuan Zhang, Limei Wang, Jacob Helwig, Youzhi Luo, Cong Fu, Yaochen Xie, Meng Liu, Yu-Ching Lin, Zhao Xu, Keqiang Yan, Keir Adams, Maurice Weiler, Xiner Li, and etc. 2023. Artificial Intelligence for Science in Quantum, Atomistic, and Continuum Systems. *ArXiv abs/2307.08423* (2023).
- [266] Yu Zhang, Xiushi Chen, Bowen Jin, Sheng Wang, Shuiwang Ji, Wei Wang, and Jiawei Han. 2024. A Comprehensive Survey of Scientific Large Language Models and Their Applications in Scientific Discovery. In *Proc. 2024 Conf. on Empir. Methods Nat. Lang. Process.*. ACL, 8783–8817.
- [267] Xueliang Zhao, Tingchen Fu, Lema Liu, Lingpeng Kong, Shuming Shi, and Rui Yan. 2023. SORTIE: Dependency-Aware Symbolic Reasoning for Logical Data-to-text Generation. In *Find. Assoc. for Comput. Linguist. ACL 2023*. ACL, 11247–11266.
- [268] Yilun Zhao, Zhenting Qi, Linyong Nan, Lorenzo Jaime Flores, and Dragomir Radev. 2023. LoFT: Enhancing Faithfulness and Diversity for Table-to-Text Generation via Logic Form Control. In *Proc. 17th Conf. Eur. Chapter Assoc. for Comput. Linguist.*. ACL, 554–561.
- [269] Yiran Zhao, Jinghan Zhang, I Chern, Siyang Gao, Pengfei Liu, Junxian He, et al. 2024. Felm: Benchmarking factuality evaluation of large language models. *Adv. Neural Inf. Process. Syst.* 36 (2024).
- [270] Yangqiaoyu Zhou, Haokun Liu, Tejes Srivastava, Hongyuan Mei, and Chenhao Tan. 2024. Hypothesis Generation with Large Language Models. *arXiv preprint arXiv:2404.04326* (2024).
- [271] Mingchen Zhuge, Changsheng Zhao, Dylan Ashley, Wenyi Wang, Dmitrii Khizbullin, Yunyang Xiong, Zechun Liu, Ernie Chang, Raghuraman Krishnamoorthi, Yuandong Tian, Yangyang Shi, Vikas Chandra, and Jürgen Schmidhuber. 2024. Agent-as-a-Judge: Evaluate Agents with Agents. *arXiv preprint arXiv:2410.10934* (2024).

Appendix

This appendix provides supplementary materials intended to support and extend the main text. It includes a background overview, further elaboration on AI support for specific topics and tasks, and a section on AI use cases that illustrates how AI tools were integrated into the workflow and phrasing of this paper.

A.1 Background

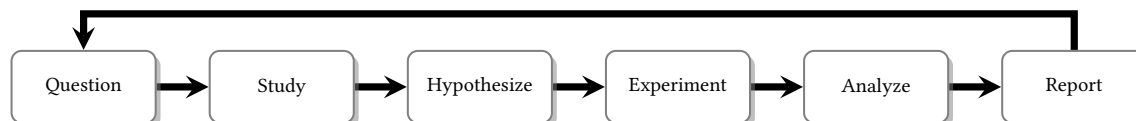


Fig. 4. Scientific discovery cycle, after [45]

Over time, science has progressed through numerous paradigm shifts, leading to the modern era of data-intensive exploration [88]. Despite advancements in tools and methodologies that have accelerated discovery, the fundamental steps of the scientific process have remained consistent. As illustrated in Fig. 4, this process typically begins with identifying a research question, followed by reviewing relevant literature, formulating a hypothesis, designing and conducting experiments, analyzing data, and ultimately reporting findings. This iterative cycle ensures the continuous refinement and expansion of scientific knowledge.

With respect to the first two of these steps, a major challenge for any scholar is achieving, and then maintaining, sufficient familiarity with existing research on a given topic to be able to identify new research questions or to discover the knowledge required to answer them. Before the 20th century, it was often feasible to keep abreast of developments in a specialty simply by reading all the relevant books and journals as they were published. In modern times, however, the number of scientific publications has been doubling every 17 years [20], making this exhaustive approach unworkable. The need to sift through large quantities of scholarly knowledge spurred the specialization of simple library catalogs (in use since ancient times) into abstracting journals, bibliographic indexes, and citation indexes. By the 1960s and 1970s, many of these resources were being produced with standardized control principles and technologies, and could be queried interactively using automated information retrieval systems [19, pp. 88–91]. These technical developments have enabled the widespread adoption of more principled approaches to the exploration of scientific knowledge, such as systematic reviews [28] and citation analysis [72].

How experts propose hypotheses to explain observed phenomena has been extensively discussed in the philosophy and psychology of science, albeit with little empirical work until relatively recently [41, 42]. Contrary to the idealized notion of scientific reasoning, hypotheses rarely come about solely through induction (i.e., the abstraction of a general principle from a set of empirical observations). Rather, case studies employing think-aloud protocols suggest that hypotheses are generated through a process of successive refinement. These processes may involve non-inductive heuristics (analogies, simplifications, imagistic reasoning, etc.) that often fail individually, but may lead to valid explanatory models after “repeated cycles of generation, evaluation, and modification or rejection” [41, 42].

Experimentation and analysis aim to establish a causal relationship between the independent and dependent variables germane to a given scientific hypothesis. The metascientific literature abounds with practical advice on the design and execution of experiments, much of it discipline-specific. However, the general ideas at play can be traced to Ronald

Fisher, whose seminal works on statistical methods [64] and experimental design [65] popularized the principles of randomization (assigning experimental subjects by chance), replication (observing different experimental subjects under the same conditions), and blocking (eliminating undesired sources of variation). Besides these considerations, experimental design involves the determination of the (statistical) analysis that will be performed, and is often constrained by the availability of resources such as the time, effort, or cost to gather and analyze observations or data [112].

The final step in the scientific cycle, reporting, encompasses the dissemination of research findings, typically but not exclusively to the wider scientific community through articles, books, and presentations. The practice of scientific communication has itself attracted scientific study, leading to descriptive and pedagogical treatments of its various processes and strategies (e.g., [83, 257]). The essential role of peer review [244] has attracted special attention, albeit more on its high-level processes, its efficacy and reliability, and its objectivity and bias rather than on how reviewers go about evaluating manuscripts and communicating this evaluation. Accordingly, technological developments in the peer review workflow have until very recently tended to focus on managing or streamlining the review process for the benefit of the editor and publisher, or on supporting open or collaborative reviewing [55, 244].

A.2 Supplement on AI Support for Specific Topics and Tasks

A.2.1 Additional Literature Search, Summarization, and Comparison

Platform		Search	Recommendations	Collections	Citation Analysis	Trending Analysis	Author Profiles	Visualization Tools	Paper Chat	Idea Generation	Paper Writing	Summarization	Paper Review	Datasets	Code Repositories	LLM Integration	Web API	Personalization	Cost	Data Source
Search Engines	Google Scholar	✓	✓	✓	✓		✓									✓	✓		Free	
	Semantic Scholar	✓	✓	✓	✓	✓	✓		✓	✓					✓	✓	✓		Free	214 million
	Baidu Scholar	✓	✓	✓	✓	✓	✓								✓	✓	✓		Freemium	680 million
	BASE	✓		✓												✓			Free	415 million
	Internet Archive Scholar	✓														✓			Free	35 million
	Scilit	✓		✓			✓												Free	172 million
	The Lens	✓		✓			✓									✓			Freemium	284 million
	Science.gov	✓						✓											Free	several million
	Academia.eu	✓		✓			✓												Freemium	55 million
	OpenAlex	✓					✓									✓			Freemium	
	AceMap	✓			✓	✓	✓	✓					✓						Free	260 million
	PubTator3	✓		✓	✓											✓			Free	6 million
Benchm.	Papers with Code	✓									✓	✓							Free	154 thousand
	ScienceAgentBench								✓		✓	✓	✓						Free	
	ORKG Benchmarks				✓	✓					✓								Free	
	Huggingface	✓		✓		✓					✓	✓							Freemium	

Table 6. Overview of popular literature search, summarization and comparison tools and their key features.

Search engines. Traditional academic search engines such as [Google Scholar](#), [Semantic Scholar](#), [Baidu Scholar](#), [Science.gov](#), and [BASE](#), as shown in Table 6, are characterized by their broad literature coverage, citation tracking capabilities, and keyword-based search functionality. Their primary advantages include extensive indexing of scholarly content, which involves aggregating and organizing vast amounts of academic documents from various sources such as publisher websites, institutional repositories, and open-access archives. This comprehensive indexing spans multiple disciplines and document types, ensuring that users can access a diverse set of resources. Additionally, these platforms

offer robust citation analysis features that allow researchers to track citation counts, measure the impact of publications, and explore citation networks to identify influential works and emerging trends within a given field. Another significant advantage is their free access to a wide range of academic resources, such as peer-reviewed journal articles, conference papers, preprints, theses and dissertations, technical reports, books and book chapters, as well as grey literature like white papers, government reports, and institutional research outputs. However, these search engines have certain limitations, such as limited AI-driven filtering options and relatively basic relevance ranking mechanisms compared to more advanced AI-enhanced search tools.

Benchmarks and leaderboards. Code and Dataset-Focused Search Engines include platforms such as [Papers with Code](#), [ScienceAgentBench](#), and [Huggingface](#), which are specifically designed to bridge the gap between academic publications and practical implementation by linking research papers with associated code and datasets. These platforms facilitate reproducibility and practical application of research findings by aggregating code repositories, enabling researchers and practitioners to easily explore implementations, compare results, and benchmark their models. A key feature of such platforms is their ability to provide dataset discovery tools, which allow users to identify relevant datasets for specific research problems, fostering collaboration and accelerating experimentation cycles. These search engines are particularly valuable for machine learning practitioners, as they facilitate quick access to ready-to-use codebases, helping them implement cutting-edge research more efficiently. Based on these community-curated leaderboards, some studies have proposed models for constructing leaderboards directly from scientific papers [91, 107, 194].

A.2.2 Extended AI-Driven Scientific Discovery: Ideation, Hypothesis Development, and Experimentation

Fig. 5 provides a broad methodological overview, grouping methods applied in hypothesis generation, idea generation, and automated experimentation. In contrast, Fig. 2 presents concrete examples and references to exemplary papers.

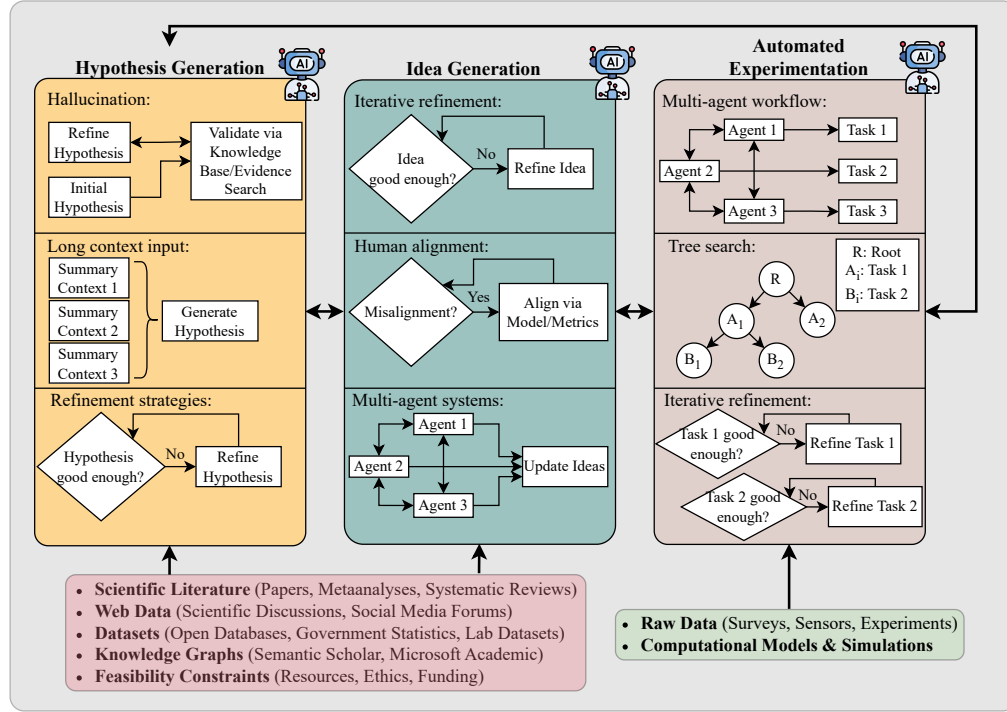


Fig. 5. Visualization of the hypothesis generation, idea generation, and automated experimentation process. Most works in hypothesis generation focus on reducing hallucinations, handling long contexts, and iteratively refining outputs. To reduce hallucinations, an initial hypothesis is validated against a knowledge base for refinement. For long-context inputs, different contexts are summarized and integrated, while refinement strategies iteratively improve the hypothesis until it meets a satisfactory level. A similar iterative refinement strategy is also applied to idea generation. Additionally, alignment strategies are employed to make generated ideas more thoughtful and feasible. In multi-agent systems, multiple agents collaborate to enhance the idea generation process. In contrast, automated experimentation often relies on tree search for selecting optimal examples, multi-agent workflows where LLMs collaborate on distinct tasks, and iterative refinement to improve task performance. While hypothesis and idea generation leverage diverse sources such as scientific literature, web data, and datasets, automated experimentation operates on predefined ideas and requires access to computational models, simulations, and raw data.

A.2.3 Text-based Content Generation

In this section, we provide an additional figure (Fig. 6) to illustrate the content generation process for academic papers, covering title, abstract, related work, and bibliography generation with their respective methods.

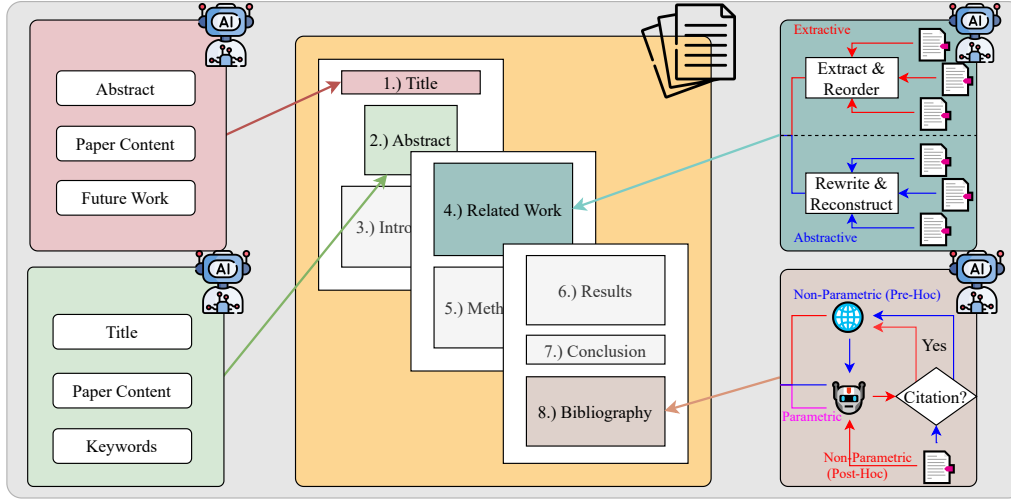


Fig. 6. Visualization of the content generation process for academic papers. Title generation methods include abstract-to-title, content-to-title, and future work-to-title mappings. Abstract generation typically involves title-to-abstract and keywords-to-abstract techniques. Related work generation follows either extractive methods (reordering extracted sentences) or abstractive methods (rewriting content from multiple papers). Bibliography generation is categorized into non-parametric methods (retrieving references from external sources) and parametric methods (LLMs generating references from preexisting knowledge without retrieval). Non-parametric methods are further divided into pre-hoc (determining citation needs before text generation and retrieving references beforehand) and post-hoc (checking for citations after text generation and appending retrieved references as needed).

A.2.4 Multimodal content generation and understanding: additional information

A.2.4.1 Methods and results. In the following, we provide a summary of representative approaches for multimodal content generation and understanding in Table 7, illustrate the process of scientific figure generation in Fig. 7, and present an extended description of scientific slide and poster generation.

Scientific slide and poster generation. For scientific slide generation, early works typically relied on heuristic rule-based approaches. For instance, Sravanthi et al. [211] develop a rule-based system to generate slides for each section and subsection of a paper. The slide text content is generated using a query-based extractive summarization system. Later, researchers began to leverage machine learning approaches to extract key phrases and their corresponding important sentences. Hu and Wan [93] use a Support Vector Regression (SVR) model to learn the importance of each sentence in a paper. The slides are then generated using an integer linear programming (ILP) model to select and align key phrases and sentences. Wang et al. [240] propose a system to generate slides for each section of a given paper, focusing on creating two-layer bullet points. The authors first extract key phrases from the paper using a parser and then use a random forest classifier to predict the hierarchical relationships between pairs of phrases. Li et al. [127] develop two sentence extractors—a neural-based model and a log-linear model—within a mutual learning framework to extract

Task	Input	Output	Dataset	Method	Evaluation
Scientific Figure Understanding					
Question Answering [108]	Synthetic, scientific-style figures and questions	Answers	FigureQA	Fine-tuning	Accuracy
Chart Summarization [185]	Chart images with metadata	Chart summaries	ChartSumm	Fine-tuning	Automatic evaluation
Caption Figure Retrieval	Figure or caption	Caption or figure	SciMMIR	Fine-tuning	Ranking Metrics
Scientific Figure Generation					
Caption/Instruction-to-code generation [14, 229]	(Extended) scientific caption or instruction	Compilable (TikZ, Vega, etc.) code of scientific figure	AutomaTikZ, DaTikZ	Fine-tuning	Human & various metrics [14]
Description-to-image generation [263]	Description/instruction	Scientific image	SciImage	Prompting	Human
Sketch/Image-to-image generation	Scientific (raster) image or sketch	Compilable TikZ code of scientific figure	DaTikZ-v2	Fine-tuning & MCTS	Human & various metrics
Scientific diagram generation [161]	Scientific paper(s) + intent	Diagram	SciDoc2-DiagramBench	Two-stage pipeline	Human & various metrics
Scientific Table Understanding					
Table description [163]	Tables from scientific articles	Table description	SciGen	Fine-tuning	Automatic & human evaluation
Numerical reasoning [216]	Tables from scientific papers	Numerical descriptions	NumericNLG	Fine-tuning	Automatic & human evaluation
Scientific Table Generation					
Literature review table generation [167]	A list of papers	Table schema + values	ArXivDigest Tables	Prompting	Automatic & human evaluation
Scientific Slide and Poster Generation					
Single slide generation [217]	Paper + slide title	Slide content	SciDuet	Two-step method	ROUGE and human evaluation
Slide deck generation [68]	Paper	A deck of slides	DOC2PPT	Hierarchical generative model	Automatic & human evaluation
Personalized slide deck generation [162]	Paper + target audience (technical or non-technical)	A deck of slides	Persona-Aware-D2S	Fine-tuning	Automatic & human evaluation

Table 7. Multimodal Content Generation and Understanding Approaches.

relevant sentences from papers. These sentences are used to generate draft slides for four topics: *Contribution*, *Dataset*, *Baseline*, and *Future Work*.

It is important to note that all the aforementioned works focus on extracting sentences or phrases from the given paper to serve as the slide text content. In contrast, Fu et al. [68] and Sun et al. [217] take a different approach by training sequence-to-sequence models to generate sentences for the slide text content. This distinction is analogous to the difference between “extractive” and “abstractive” summaries in text summarization. More specifically, Fu et al. [68] design a hierarchical recurrent sequence-to-sequence architecture to encode the input document, including sentences

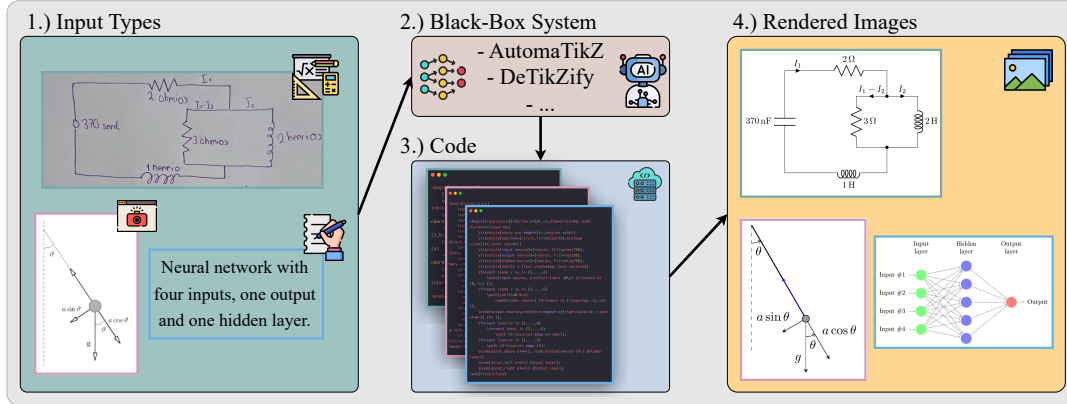


Fig. 7. Overview of the scientific figure generation process. Various input types including sketches, screenshots, and text, can be used to generate TikZ code with tools such as AutomaTikZ [14] and DeTikZify [15]. The generated code is then rendered into high-quality vector graphics images.

and images, and generate a slide deck. In contrast, Sun et al. [217] assume that slide titles would be provided by end users. The authors use these titles to retrieve relevant and engaging text, figures, and tables from the given paper using a dense retrieval model. They then summarize the retrieved content into bullet points with a fine-tuned long-form question answering system based on BART.

With recent advancements in LLMs and vision-language models (VLMs), researchers have started utilizing these technologies for generating scientific presentation slides. Mondal et al. [162] propose a system to generate persona-aware presentation slides by fine-tuning LLMs such as *text-davinci-003* and *gpt-3.5-turbo* with a small training dataset containing personalized slide decks for each paper. Maheshwari et al. [158] focus solely on generating text content and develop an approach that combines graph neural networks (GNNs) with LLMs to capture non-linearity in presentation generation, while attributing source paragraphs to each generated slide within the presentation. Bandyopadhyay et al. [11] design a bird's-eye view document representation to generate an outline, map slides to sections, and then create text content for each slide individually using LLMs. The approach then extracts images from the original papers by identifying text-image similarity in a shared subspace through a VLM.

Generating posters from scientific papers has received less attention compared to scientific slide generation. Qiang et al. [183] introduce a graphical model to infer key content, panel layouts, and the attributes of each panel from data. Xu and Wan [250] develop a demo system for automated poster generation. The system first identifies important sections using a trained classifier. It then employs a summarization model to extract key sentences and related graphs from each section to construct corresponding panels. Finally, the system generates a LaTeX document for the poster based on the template selected by the user.

A.3 AI use case and abbreviations

Throughout this paper, we have integrated AI tools to support specific aspects of the research workflow. For example, in the subsection of *Literature Search, Summarization, and Comparison*, we used Google Search, ChatGPT, NotebookLM, and Scholar Inbox to retrieve relevant tools and related work. Additionally, LLMs assisted with grammar and spell checking, as well as generating code for formatting tables.

The abbreviations used in our paper are summarized in Table 8.

Acronym	Full Name
AI4Research	Towards a Knowledge-grounded Scientific Research Lifecycle
AISD	AI & Scientific Discovery
CoI	Chain of Ideas
CV	Computer Vision
FM4Science	Foundation Models for Science
GNNs	Graph Neural Networks
ILP	Integer Linear Programming
KGs	Knowledge Graphs
LLMs	Large Language Models
LSTM	Long Short-Term Memory Networks
ML	Machine Learning
MSE	Mean-Square Error
MTEB	Massive Text Embedding Benchmark
NLP	Natural Language Processing
NSLP	Natural Scientific Language Processing and Research Knowledge Graphs
QA	Question Answering
RAG	Retrieval-Augmented Generation
S2ORC	Scholar Open Research Corpus
SVR	Support Vector Regression
TFR	Text-Figure Relevance
VLMs	Vision-Language Models

Table 8. List of Acronyms and Their Full Names