

Position: Multimodal Large Language Models Can Significantly Advance Scientific Reasoning

Anonymous ACL submission

Abstract

Scientific reasoning, the process through which humans apply logic, evidence, and critical thinking to explore and interpret scientific phenomena, is essential in advancing knowledge reasoning across diverse fields. However, despite significant progress, current scientific reasoning models still *struggle with generalization across domains and often fall short of multimodal perception*. Multimodal Large Language Models (MLLMs), which integrate text, images, and other modalities, present an exciting opportunity to overcome these limitations and enhance scientific reasoning. Therefore, **this position paper argues that MLLMs can significantly advance scientific reasoning** across disciplines such as mathematics, physics, chemistry, and biology. First, we propose a four-stage research roadmap of scientific reasoning capabilities, and highlight the current state of MLLM applications in scientific reasoning, noting their ability to integrate and reason over diverse data types. Second, we summarize the key challenges that remain obstacles to achieving MLLM’s full potential. To address these challenges, we propose actionable insights and suggestions for the future. Overall, our work offers a novel perspective on MLLM integration with scientific reasoning, providing the LLM community with a valuable vision for achieving Artificial General Intelligence (AGI).

1 Introduction

Scientific reasoning, at its core, is the process through which humans apply logic, evidence, and critical thinking to explore and interpret phenomena in various scientific domains (Bao et al., 2009; Lawson, 2004). This cognitive ability is essential not only for advancing knowledge but also for fostering a deeper understanding of the natural world, particularly in fields such as mathematics, physics, chemistry, and biology. In education, medicine,

finance, AI for Science, and other domains, scientific reasoning serves as a cornerstone for cultivating problem-solving skills, analytical thinking, and innovation (Jadon et al., 2025). However, despite shared objectives, each domain has unique characteristics in terms of data representation, knowledge construction, and reasoning methods (Ferrag et al., 2025; Chen et al., 2025).

In response to these challenges, the scientific community has explored a range of approaches, from traditional statistical methods to the more recent advancements in deep learning, with the goal of improving knowledge reasoning across disciplines (Goodman, 2016; Lu et al., 2022b). While significant progress has been made in enhancing scientific reasoning within specific domains, a gap remains in the broader context of scientific research. Current scientific reasoning models, and even those targeted toward domain-specific applications, are still *far from achieving the generalization capabilities necessary for Artificial General Intelligence (AGI), which aims to exhibit unified reasoning across all fields* (Birhane et al., 2023).

The rapid rise of Large Language Models (LLMs) in recent years has brought transformative changes across various domains, pushing the boundaries of what is possible in natural language processing and understanding (Min et al., 2023; Zhao et al., 2023a). Despite their remarkable zero-shot reasoning abilities, many areas, particularly in scientific fields, require multimodal inputs to build a comprehensive understanding of knowledge. This has led to the emergence and growth of Multimodal Large Language Models (MLLMs), which are capable of integrating and reasoning over multiple types of data, such as text, images, and other modalities (Liang et al., 2024b; Bai et al., 2024a). MLLMs are not only revolutionizing language understanding but also paving the way for advancements in scientific reasoning by processing complex multimodal data in ways that were

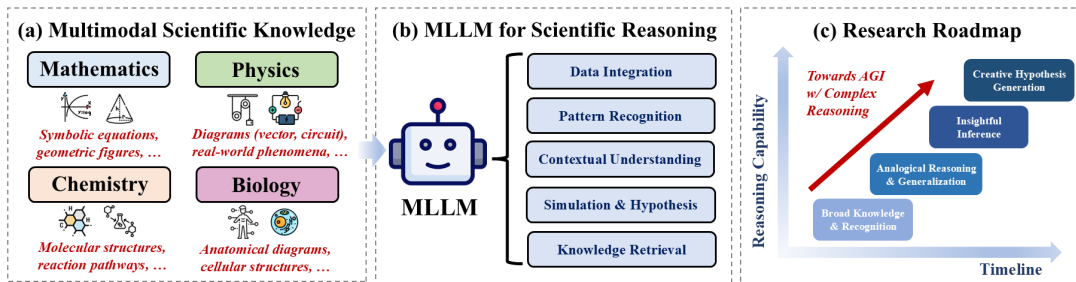


Figure 1: The big picture of our position. We focus on multimodal scientific fields, especially mathematics, physics, chemistry, and biology as our scope (a), and we advocate leveraging MLLMs with multiple reasoning functions for scientific reasoning (b). We further propose a four-stage roadmap for scientific reasoning capability, ultimately achieving AGI (c).

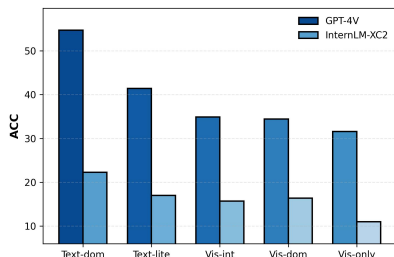


Figure 2: MLLM performance declines w/ increased reliance on visual modality from MathVerse (Zhang et al., 2025).

previously unachievable. However, a critical gap in visual reasoning capabilities persists, as MLLM performance degrades when shifting reliance from textual descriptions to visual diagrams (Figure 2).

Building on the rapid evolution of MLLMs and the growing demand for enhanced multimodal reasoning capabilities in scientific domains (Reddy and Shojaee, 2024; Zhang et al., 2024d), **this paper proposes a position: MLLMs can significantly advance scientific reasoning.** By integrating multimodal learning techniques, MLLMs have the potential to address the pressing challenges in scientific reasoning, as shown in Figure 1. This paper aims to break down the analysis into three key sections: the current state of MLLMs in scientific reasoning, the challenges encountered, and the future steps needed to achieve greater success.

Our analysis delves into these topics in detail, starting with an overview of background and research roadmap (Section 2), followed by how MLLMs are currently handling scientific reasoning (Section 3) as well as a discussion of the challenges faced (Section 4), and concluding with the promising opportunities that lie ahead (Section 5). In particular, we first explore how MLLMs have been applied across scientific disciplines, detailing the techniques used for data integration, pattern recognition, contextual understanding, and more. The Challenges section examines the inherent difficulties, ranging from technical to ethical aspects, that need to be overcome for MLLMs to achieve

optimal performance. Finally, in the Discussion section, we highlight future steps with relevant literature, from data & training strategies to agent-based collaboration, as key areas for advancing the integration of MLLMs with scientific reasoning. Through this, we aim to provide a comprehensive understanding of the landscape and offer strategic insights for future research in the intersection of MLLMs and scientific reasoning¹.

Our position and analytical framework contribute in three significant ways. ❶ We present a **novel perspective** on the integration of MLLMs with scientific reasoning, highlighting how this synthesis could reshape research in the field. ❷ We offer a **systematic review and categorization** of recent advancements, showcasing representative work and outlining a clear roadmap for future progress. ❸ We identify and explore **future opportunities**, providing the community with promising insights that could guide the next generation of research in scientific reasoning.

2 Background

2.1 Scientific Reasoning

Scientific reasoning is the intellectual process of forming hypotheses, interpreting evidence, and applying logical frameworks to solve problems or explain phenomena (Bao et al., 2009; Lawson, 2004). Its importance spans diverse scientific domains, such as mathematics, physics, chemistry, and biology, where it drives discovery, fosters understanding, and enables practical innovation. With the rise of multimodal data, scientific reasoning increasingly requires integrating and synthesizing information from multiple sources, including textual, visual, and other modalities².

The significance of scientific reasoning in the MLLM context is profound. By enabling mod-

¹See more clarification of our scope in Appendix A.

²See the formal formulation of the task in Appendix B.

els to connect disparate data points and infer relationships across modalities, MLLMs hold the potential to transform how researchers approach interdisciplinary problems. This capability is critical for addressing grand challenges such as climate modeling, drug discovery, *etc* (Zhang et al., 2024d). Moreover, enhancing MLLM-based scientific reasoning aligns with broader goal of advancing AGI, as it exemplifies synthesis of learning, abstraction, and decision-making across domains (Guo et al., 2025a; Thawakar et al., 2025).

2.2 Multimodal Large Language Models

Most existing MLLMs consist of three primary modules: a modality encoder, an LLM module, and a projector between them (Fu et al., 2024). Typically, the modality encoder extracts embeddings from non-language modalities such as images or audio, which are then projected into the word space of the LLM via the projector. The post-projection embeddings are subsequently combined with word embeddings derived from system prompts and user queries to serve as input for the LLM. Similar to LLMs, MLLMs generate responses in an autoregressive manner:

$$p(\mathbf{w}_O \mid \mathbf{w}_V, \mathbf{w}_T) \sim \prod_{t=1}^L P(w_t \mid w_{<t}, \mathbf{w}_V, \mathbf{w}_T)$$

Here, \mathbf{w}_V and \mathbf{w}_T denote the post-projection embeddings and word embeddings respectively, while $\mathbf{w}_O = \{w_{o,t}\}_{t=1}^L$ represents the generated word token sequence of length L . With their capability to comprehend visual inputs, contemporary MLLMs demonstrate remarkable performance in various tasks, including visual question answering (VQA) (Ishmam et al., 2024; Uppal et al., 2022; Dang et al., 2024a), image captioning (Vaishnavi and Narmatha, 2024; Agarwal and Verma, 2024), and multimodal reasoning (Yan and Lee, 2024; Yan et al., 2024a; Huo et al., 2024).

Presently, there is a plethora of open-source foundation MLLMs capable of general multimodal tasks. Notable examples include LLaVA family (Liu et al., 2023b,a, 2024b), Qwen-VL series (Bai et al., 2023; Wang et al., 2024e), InternVL series (Chen et al., 2024b,a), LLaMA-3.2-Vision (Dubey et al., 2024), *etc*. Despite these advancements, open-source MLLMs still lag behind closed-source models like GPT-4o (Achiam et al., 2024), Claude (Anthropic, 2024), and Gemini-Pro (Team et al., 2024a) in complex reasoning

tasks (Liu et al., 2025d; Yue et al., 2024a). With the emergence of o1-like reasoning models (Jaech et al., 2024; Zeng et al., 2025; Liu et al., 2025b), preliminary efforts are underway to elicit the slow-thinking capabilities of MLLMs, as seen in works like QvQ (Qwen, 2024a), Mulberry (Yao et al., 2024), Virgo (Du et al., 2025), *etc*.

3 How MLLMs Benefit Scientific Reasoning

3.1 Research Roadmap

The development of (M)LLMs for scientific reasoning can be categorized into four progressive stages: *Broad Knowledge and Recognition*, *Analogical Reasoning and Generalization*, *Insightful Inference*, and *Creative Hypothesis Generation*. Each stage is defined by its unique characteristics across four dimensions: data and knowledge requirements, reasoning mechanisms, model generalization, and applications and impact (See the detailed comparisons among the four dimensions in Appendix C). Figure 1 (c) provides an overview of four stages, highlighting their evolution progress.

Stage 1: Broad Knowledge and Recognition. The initial stage focuses on building a strong foundational understanding across domains. MLLMs in this stage rely on highly diverse and multimodal datasets to capture a broad range of knowledge. Reasoning mechanisms are primarily retrieval-based, with emphasis on pattern recognition, data alignment, and summarization. Model generalization remains limited, operating primarily within predefined domains (White, 2023; Pei et al., 2024; Chen et al., 2023c).

Stage 2: Analogical Reasoning and Generalization. This stage emphasizes the ability to draw connections and analogies across domains. Data requirements shift towards moderately diverse datasets that emphasize relationships and cross-domain patterns. Reasoning mechanisms incorporate relational reasoning and analogical thinking, enabling MLLMs to generalize effectively across domains. Applications include interdisciplinary problem-solving, transfer learning, and identifying cross-domain insight, reflecting a moderate increase in complexity and impact (Webb et al., 2023; Lewis and Mitchell, 2024).

Stage 3: Insightful Inference. The third stage focuses on inferring deep insights from minimal and high-context data. Data requirements narrow to low-diversity, domain-specific datasets, al-

lowing MLLMs to develop nuanced understanding. Reasoning mechanisms involve predictive reasoning and contextual interpretation, enabling the model to deduce complex outcomes. Generalization becomes highly context-specific, and applications include optimization and predictive modeling, making this stage highly impactful (Melko and Carrasquilla, 2024; Barman et al., 2025).

Stage 4: Creative Hypothesis Generation. In the final stage, MLLMs achieve the ability to generate innovative hypotheses and explore uncharted territories. Data requirements include highly diverse and synthetic datasets or simulation environments, fostering creativity. Reasoning mechanisms reach their highest complexity, involving generative reasoning and hypothesis exploration. Generalization becomes innovation-driven, synthesizing knowledge across fields. Applications at this stage have the highest impact, including proposing new theories, designing experiments, and driving scientific discovery (Xiong et al., 2024a; Qi et al., 2024; Pelletier et al., 2024).

3.2 Data Heterogeneity Across Four Scientific Domains

MLLMs are designed to process and integrate information from both textual and visual modalities, offering a versatile framework for handling complex scientific reasoning. However, the distinct nature of data across disciplines introduces unique challenges in model training and application. Appendix D summarizes key differences in visual features for four scientific subjects within our scope.

Each subject presents unique challenges in data representation, with mathematical tasks primarily focusing on abstract symbols and formulas, while other subjects, particularly biology, require a mix of detailed real-world imagery and conceptual explanations. These differences necessitate domain-specific adaptations in how MLLMs process and understand multimodal data (Bai et al., 2024a).

3.3 MLLM-based Scientific Reasoning

As shown in Figure 3, current MLLM-based scientific reasoning can generally be divided into the following five paradigms, which progressively enhance the reasoning capabilities of MLLMs, ultimately moving towards AGI.

Data Integration. One of the primary strengths of MLLMs is their ability to integrate multimodal data from various scientific domains. For example, in physics, models can combine textual de-

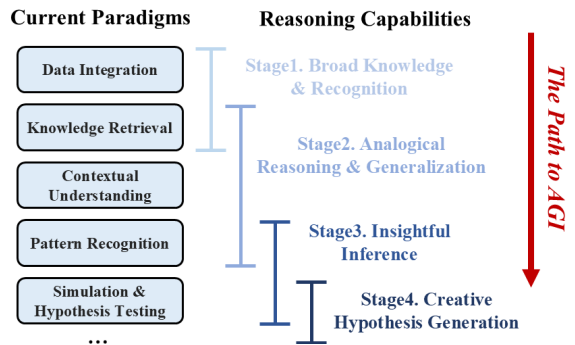


Figure 3: Overview of MLLM-based scientific reasoning paradigms and corresponding reasoning capabilities.

scriptions of a problem, such as Newtons second law, with visual representations like force diagrams. The model can then reason about how different forces interact and predict motion (Barman et al., 2025; Sato, 2024). Similarly, in chemistry, MLLMs can combine chemical equations with 3D molecular structures, offering deeper insights into reaction mechanisms (Guo et al., 2023). This integrated approach allows MLLMs to handle intricate, often disjointed, data sources to generate coherent scientific explanations.

The ability to integrate and synthesize multimodal information enables the MLLM to solve complex problems more effectively. However, challenges arise when the visual data does not perfectly align with the textual explanation, potentially leading to misinterpretations (Zhang et al., 2025; Lu et al., 2023; Zhuang et al., 2024).

Knowledge Retrieval. A significant aspect of scientific reasoning is knowledge retrieval. In fields like physics and chemistry, the vast amount of scientific knowledge available - such as established theories, laws, or empirical data - can be overwhelming. Knowledge retrieval helps MLLMs access external knowledge bases, databases, and scientific literature to supplement their reasoning. For instance, when solving a chemistry problem, an MLLM could retrieve data from a chemical database to identify missing properties of substances or reactions that were not explicitly stated in the task (Prince et al., 2024; Sze and Hassoun, 2024). In addition, knowledge retrieval can aid MLLMs in bridging gaps between modalities. For example, an MLLM working on a biological problem may retrieve relevant studies from scientific papers to fill in missing knowledge, such as identifying unknown interactions between proteins or cells (Li et al., 2024j, 2025).

This aspect of MLLMs reasoning ensures the model remains up-to-date with the latest discov-

eries and can apply a deeper layer of scientific knowledge in reasoning processes. However, challenges in accurately selecting and integrating relevant knowledge remain, particularly when sources of conflict are present (Fan et al., 2024).

Contextual Understanding. Contextual understanding in scientific reasoning involves understanding not only the literal data presented but also the broader context in which it is used. MLLMs are capable of this by combining visual data, such as molecular structures in chemistry, with textual descriptions of chemical properties. This allows them to reason about potential interactions between molecules in a way that goes beyond simple matching (Liu et al., 2025a; Horawalavithana et al., 2023). Wang et al. (2024b) leverage Chain-of-Thought as teaching signals to train small models to perform reasoning in complicated scenarios.

This contextual capability is crucial in fields like biology, where visual images of biological processes must be linked with underlying theories to make accurate predictions. However, this capability can be limited when the model fails to integrate textual and visual information effectively, leading to errors in reasoning (Li and Tang, 2024).

Pattern Recognition. In scientific reasoning, pattern recognition is a crucial skill that MLLMs excel at. MLLMs can detect patterns across different modalities, whether they are geometric in mathematics or experimental in chemistry. For instance, in biology, MLLMs can recognize cellular structures in images and relate them to known biological processes described in text, such as identifying mitochondria and correlating their function with energy production (Luu and Buehler, 2023; Kraus et al., 2024). Additionally, an example of pattern recognition in mathematics could involve an MLLM matching visual representations of geometric figures with algebraic equations to find solutions to geometry problems (Mouselinos et al., 2024). This capability enhances model’s ability to understand complex systems across disciplines, including identifying patterns in large datasets that might be intricate for manual analysis.

This skill allows MLLMs to be highly effective in analyzing multistep reasoning problems, which are often required in scientific disciplines (Qiao et al., 2024; Yan et al., 2024b). However, pattern recognition can be hindered by noisy or low-quality visual data, particularly in domains like biology, where image clarity is critical for correct interpretation (Zhang et al., 2024e; Ren et al., 2024).

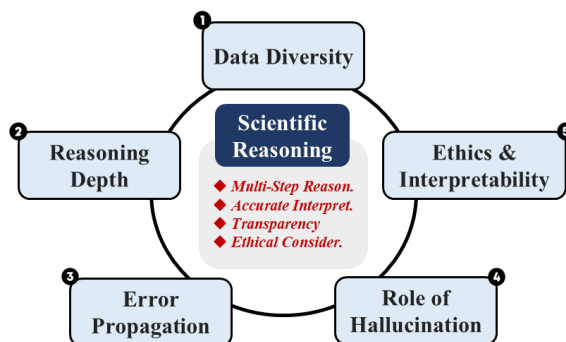


Figure 4: Challenges for MLLM-based scientific reasoning.

Simulation and Hypothesis Testing. MLLMs also possess the ability to perform simulation and hypothesis testing, a fundamental part of scientific reasoning (Qi et al., 2023). For example, in physics, MLLMs can simulate the effect of various forces on an object, predict outcomes, and validate those predictions against real-world data or experiments (Melko and Carrasquilla, 2024; Gao et al., 2024a; Yu et al., 2024). This capacity allows MLLMs to conduct scientific inquiries in a manner akin to human researchers, testing hypotheses and refining conclusions (Morera, 2024).

Despite these strengths, hypothesis testing is constrained by the quality & quantity of available data, which in some scientific domains remains insufficient or incomplete. This limits the generalizability and reliability of MLLMs in tasks requiring deep, multistep reasoning (Xiong et al., 2024a).

4 Challenges

Even though MLLMs show substantial promise in solving scientific reasoning tasks, significant challenges remain. Based on the inherent characteristics of scientific reasoning - the need for multi-step inference and precise speculation, while ensuring transparency and ethicality - we further propose the following five key challenges (Figure 4).

Data Diversity. Another challenge is the diversity of data across different scientific domains (summary of multi-domain scientific datasets can be seen in Appendix E). While mathematics is rich in textual data, such as equations and proofs, the availability of high-quality visual data is more limited (Qiao et al., 2024; Sun et al., 2024a; Liu et al., 2024c; He et al., 2024a). In contrast, fields like chemistry and biology benefit from abundant visual data, such as molecular structures and microscopic images, but the corresponding textual descriptions may not always provide the depth required for comprehensive reasoning (Alampara

et al., 2024; Hocky, 2024). Without sufficient high-quality data for all modalities, model’s ability to generalize across domains is compromised.

Reasoning Depth. MLLMs frequently struggle with tasks that require deep, multi-step reasoning, especially when abstract concepts are involved. In mathematics, for example, solving a theorem involves a series of logical steps that must be followed precisely (Chen et al., 2023b). In physics, simulating complex systems, such as thermodynamics or quantum mechanics, requires a deep understanding of abstract principles and real-world conditions (Melko and Carrasquilla, 2024). MLLMs often fail to maintain this depth of reasoning, especially when applied to tasks involving intricate concepts or lengthy proof processes. This issue is particularly prevalent in fields where the complexity of reasoning extends beyond surface-level analysis and requires models to maintain rigorous logical consistency. Therefore, recent work has focused on two directions to improve the length of correct reasoning chains: the development of high-quality reasoning process datasets (Yan et al., 2024b) and the introduction of process reward models (o1 Team, 2024).

Error Propagation. Error propagation is another significant challenge in multimodal reasoning. Errors in one modality, such as a misinterpreted graph or an unclear image, can propagate throughout the reasoning process, leading to incorrect conclusions (Li et al., 2024i; Yan et al., 2024b; Li et al., 2024b). For example, in a physics problem involving force vectors, an error in interpreting the vector diagram could lead to an incorrect calculation of the net force, which would then affect subsequent steps in the solution process (Jaiswal et al., 2024). The risk of error propagation is particularly high when models are tasked with handling complex, multistep problems across multiple modalities. In particular, the impact of error propagation is especially acute in fields like physics and chemistry, where the accuracy of one step can influence the entire solution process. Small errors in initial data interpretation can lead to significant discrepancies in the final outcome (Xu et al., 2024a; Li et al., 2024c).

Role of Hallucinations. One of the most complex challenges in leveraging MLLMs for scientific reasoning is determining whether hallucinations, the generation of information not grounded in the input data or knowledge base, are inherently harmful or potentially beneficial (Bai et al.,

2024b; Liu et al., 2024a). While hallucinations are widely regarded as detrimental in factual tasks, their role in scientific reasoning is nuanced, particularly when considering the ultimate goal of advancing to Stage 4 of the research roadmap (*i.e.*, Creative Hypothesis Generation) (Jiang et al., 2024a). In scientific reasoning, hallucinations can undermine trust and reliability by introducing inaccuracies in critical domains. For example, in physics, a hallucinated formula or principle might lead to invalid conclusions, while in chemistry, a fabricated reaction pathway could suggest impossible or even dangerous experiments. These inaccuracies not only hinder immediate problem-solving but also propagate errors if used as a basis for further research (Li et al., 2024a; Chakraborty et al., 2024; Xu et al., 2024c). See details of hallucinations in scientific reasoning in Appendix F.

Ethical and Interpretability Issues. Ethical concerns and model interpretability are major challenges when deploying MLLMs in high-stakes scientific domains, such as medical research or chemical engineering. MLLMs often lack transparency, making it difficult for users to understand how the model arrived at a particular conclusion (AlSaad et al., 2024). Furthermore, ethical concerns arise when MLLMs are used to make decisions that could have significant consequences, such as in medical diagnoses or environmental impact assessments (Rahman et al., 2024). In biology and medicine, the potential for biased reasoning in MLLMs, especially when trained on unbalanced datasets, could lead to harmful or misleading conclusions (Stureborg et al., 2024; Wang et al., 2024d). An MLLM trained on biased medical data could fail to recognize critical symptoms in underrepresented populations, leading to erroneous diagnoses or treatment recommendations.

5 Discussion: What Next?

Building on the challenges outlined in Section 4, it is evident that while MLLMs hold great promise in advancing scientific reasoning, targeted solutions must be developed to address the limitations. Thus, we explore eight key perspectives for improving MLLMs in scientific reasoning (Figure 5).

The Necessity of Unified Scientific MLLMs. Although many domain-specific models have achieved remarkable performance in specialized scientific fields, exploring unified scientific MLLMs remains a critical pursuit (Taylor et al., 2022). Domain-specific models are optimized for

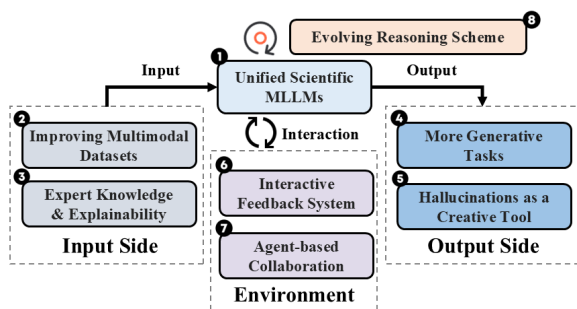


Figure 5: Eight prospects for the future of MLLMs in the field of multimodal scientific reasoning.

particular areas (summary of scientific MLLMs can be seen in Appendix G), but they often lack the ability to integrate knowledge across disciplines (Wang et al., 2023b; Shi et al., 2024a). In contrast, a unified MLLM could facilitate interdisciplinary reasoning, leveraging connections between fields to tackle complex problems that require holistic understanding, *e.g.*, climate change modeling (Nguyen et al., 2023) or biomedical research (Wang et al., 2023a).

For example, a unified scientific MLLM could simultaneously analyze chemical reaction pathways and their biological implications, enabling breakthroughs in drug discovery (Oniani et al., 2024; Guan and Wang, 2024). Similarly, it could integrate physics-based simulations with mathematical optimization to design more efficient renewable energy systems (Gao et al., 2024b).

Improving Multimodal Datasets. A critical step in advancing MLLMs’ capabilities is the improvement of multimodal datasets (Gadre et al., 2024; Rahate et al., 2022; Bayouhd et al., 2022). Current datasets often lack the richness and variety required to train models effectively across disciplines. For instance, in chemistry, existing datasets may focus heavily on molecular structures without providing sufficient textual descriptions of reaction mechanisms (Cao et al., 2023a). Similarly, in biology, while there are abundant images of anatomical structures, these are often not paired with detailed descriptions of biological processes (Tang et al., 2023; Zhang et al., 2024g). By creating multimodal datasets that include both high-quality images and comprehensive textual descriptions across all domains, the model can better learn to correlate visual features with corresponding scientific explanations (Albalak et al., 2024; Muenighoff et al., 2023). See discussion of integration with more modalities in Appendix H.

For example, in physics, datasets that combine experimental setups with corresponding theoretic

cal explanations could help improve a models understanding of underlying principles, such as energy conservation or force interactions. This integration would facilitate a more robust training process, ensuring that models can handle a wider range of scientific reasoning tasks. As a concrete suggestion, developing domain-specific multimodal datasets that cover not only the subject but different teaching contexts (*e.g.*, beginner vs. advanced materials) would help MLLMs generalize across varying complexity levels (Shi et al., 2023).

Integrating Expert Knowledge and Explainability. Integrating expert knowledge into MLLMs can significantly enhance their ability to reason accurately and logically (Pan et al., 2024). Expert knowledge, such as specialized theories in physics or established principles in biology, provides a framework within which the MLLM can operate. Additionally, by integrating causal reasoning into MLLMs, the models can better explain the relationships between variables, such as causes and effects, rather than simply identifying correlations (Xiong et al., 2024b; Jin et al., 2024b).

For example, in chemistry, integrating knowledge about chemical bonding, molecular dynamics, and reaction mechanisms can help guide the models predictions and interpretations (Zaki et al., 2024). Additionally, enhancing the explainability of MLLMs is crucial for transparency (Dang et al., 2024b). Ensuring that the model can justify its reasoning in a way that humans can easily understand, especially when making scientific claims, will increase its credibility and trustworthiness.

Expanding Scientific Reasoning to Generative Tasks. While significant progress has been made in using MLLMs for problem-solving, error detection, and theorem proving (Yan et al., 2024a), the potential of these models in generative tasks remains underexplored. Generative tasks such as creating curriculum-aligned questions (Mulla and Gharpure, 2023), designing comprehensive syllabi (Hu et al., 2024), or enabling digital teaching assistants (Onu et al., 2024) are highly relevant to real-world applications, *esp.* in educational contexts.

For example, MLLMs could be used to generate topic-specific exam questions that align with diverse educational standards or assist in creating adaptive learning materials tailored to different proficiency levels (Denny et al., 2024).

Hallucinations as a Creative Tool. Hallucinations may play a constructive role in fostering creativity and innovation, particularly in ex-

614 ploratory scientific tasks (Huang et al., 2023; Li
615 et al., 2024f). At the frontier of Stage 4 (Sec-
616 tion 3.1), MLLMs can hypothesize beyond exist-
617 ing knowledge, where generating plausible but un-
618 verified information might stimulate novel ideas.

619 A critical challenge lies in striking the balance
620 between mitigating harmful hallucinations and
621 leveraging beneficial ones. This requires a fine-
622 grained approach to model design, where halluci-
623 nation mechanisms are carefully controlled based
624 on task context. For foundational reasoning tasks,
625 strict adherence to validated knowledge is essen-
626 tial, while exploratory tasks allow controlled devi-
627 ations to inspire innovation (Jiang et al., 2024a).

628 **Interactive Feedback Systems.** Another power-
629 ful strategy involves the development of inter-
630 active feedback systems, which would allow
631 MLLMs to engage with users dynamically, iterat-
632 ing on their answers based on user input or feed-
633 back (Abramson et al., 2022; Shtarbanov et al.,
634 2023). This interactive feature would not only
635 enable models to adjust their reasoning during
636 the problem-solving process but also allow them
637 to ask clarifying questions, improving the overall
638 user experience and enhancing the models output.

639 For example, in biology, a researcher could in-
640 put a biological query related to an ecological
641 model and receive initial results from the model.
642 This back-and-forth interaction would provide a
643 mechanism for error correction and refinement, en-
644 suring that outputs align more closely with expert
645 understanding (Pan et al., 2023).

646 **Agent-based Collaboration.** A promising av-
647 enue for future development is agent-based collab-
648 oration, which involves the integration of multi-
649 ple specialized agents working together to solve
650 complex scientific problems (Guo et al., 2024b;
651 Wang et al., 2024c). Each agent could be tailored
652 to handle specific scientific tasks, such as math-
653 ematical reasoning, chemical reaction prediction,
654 or biological system analysis. These agents could
655 communicate and collaborate with each other to
656 cross-check information, validate hypotheses, and
657 combine knowledge from their respective domains
658 (Chen et al., 2023a; Yan et al., 2025).

659 For instance, in a physics problem involving
660 both mechanics and electromagnetism, an agent
661 focused on classical mechanics could collaborate
662 with an agent specialized in electromagnetism to
663 deliver a comprehensive solution that accounts
664 for the interactions between mechanical forces
665 and electromagnetic fields. By building a system

666 where multiple agents, each bringing a unique ex-
667 pertise, work collaboratively, MLLMs could ap-
668 proach complex scientific reasoning tasks with
669 higher accuracy and robustness (Guo et al., 2024b;
670 Zhao et al., 2024a).

671 **Evolving Reasoning Schemes.** Current reason-
672 ing architectures in MLLMs remain constrained
673 by opaque, monolithic designs that limit adapt-
674 ability to diverse scientific domains (Besta et al.,
675 2025). Existing paradigms of Reasoning Large
676 Models (RLMs) bifurcate into implicit RLMs (e.g.,
677 QwQ (Qwen, 2024b)), where reasoning is embed-
678 ded in model weights as a black box, and explicit
679 RLMs (e.g., LLaMA-Berry (Zhang et al., 2024b)
680 & o1 (Jaech et al., 2024)), which deploy structured
681 reasoning strategies like Monte Carlo Tree Search
682 (MCTS) or Beam Search with modular compo-
683 nents. While explicit methods enable stepwise
684 evaluation, their reliance on fixed templates and
685 proprietary training schemes hinders reproducibil-
686 ity and domain-specific customization.

687 To advance scientific reasoning, future MLLMs
688 should integrate three innovations: (i) dynamic
689 reasoning structures (e.g., nested graphs) that
690 adapt to multimodal inputs; (ii) process-based su-
691 pervision with stepwise uncertainty metrics (e.g.,
692 token-level entropy) to refine domain-specific rea-
693 soning paths; and (iii) open-source, composable
694 toolkits for hybrid training (i.e., Supervised Fine-
695 Tuning + Reinforcement Learning phases) that de-
696 couple policy/value models, enabling collabora-
697 tive, cost-efficient optimization across scientific
698 disciplines (Besta et al., 2025).

6 Alternative Views 699

700 We also discuss two key opposing perspectives
701 and provide counterarguments to address these
702 concerns in Appendix I due to the space limit.

7 Conclusion 703

704 This position paper aims to emphasize the trans-
705 formative potential of MLLMs in advancing sci-
706 entific reasoning across diverse domains, includ-
707 ing mathematics, physics, chemistry, and biology.
708 Our key position is that MLLMs represent a sig-
709 nificant step forward in enabling more comprehen-
710 sive and accurate reasoning about scientific phe-
711 nomena, bridging gaps between different types
712 of data and reasoning methods. To support this
713 stance, we reviewed the current state of MLLM
714 applications in scientific reasoning, outlined key
715 challenges, and proposed actionable insights.

716 Limitations

717 While this position paper presents a comprehen-
718 sive vision for the role of MLLMs in scientific rea-
719 soning, we acknowledge several limitations that
720 define the boundaries of our current analysis and
721 highlight avenues for future exploration.

- 722 • **Focused Disciplinary Scope.** Our analysis
723 is anchored in four core scientific disciplines:
724 mathematics, physics, chemistry, and biology.
725 While these fields are foundational and highly
726 representative of multimodal reasoning chal-
727 lenges, they do not encompass the full spec-
728 trum of scientific inquiry. Future work could
729 extend our proposed framework to other do-
730 mains, such as earth sciences, materials sci-
731 ence, and the social sciences, which present
732 their own unique data modalities and reason-
733 ing paradigms.
- 734 • **Conceptual Nature of the Roadmap.** The
735 proposed four-stage research roadmap is in-
736 tended as a high-level conceptual framework
737 to guide future development. We recognize
738 that the progression between stages may not
739 be strictly linear and the boundaries can be
740 fluid. Establishing fine-grained, quantitative
741 metrics to precisely benchmark an MLLM’s
742 position within this roadmap is a complex
743 challenge that we leave for future research.
- 744 • **Focus on Model Capabilities over Human-
745 AI Interaction.** Our discussion centers pre-
746 dominantly on the intrinsic reasoning capa-
747 bilities of MLLMs. While we briefly touch
748 upon interactive systems, a deeper explo-
749 ration of the socio-technical dynamicshow
750 these advanced models will effectively and
751 ethically integrate into the daily workflows of
752 human scientists and foster optimal human-
753 AI collaborationis beyond our current scope.
754 This represents a rich and vital avenue for fu-
755 ture work at the intersection of AI, HCI, and
756 the philosophy of science.

757 References

758 Josh Abramson, Arun Ahuja, Federico Carnevale,
759 Petko Georgiev, Alex Goldin, Alden Hung, Jessica
760 Landon, Jirka Lhotka, Timothy Lillicrap, Alistair
761 Muldal, and 1 others. 2022. Improving multimodal
762 interactive agents with reinforcement learning from
763 human feedback. *arXiv preprint arXiv:2211.11602*.

764 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama
765 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
766 Diogo Almeida, Janko Altschmidt, Sam Altman,
767 Shyamal Anadkat, and 1 others. 2024. Gpt-4 techni-
768 cal report. *arXiv preprint arXiv:2303.08774*.

769 Lakshita Agarwal and Bindu Verma. 2024. From
770 methods to datasets: A survey on image-caption
771 generators. *Multimedia Tools and Applications*,
772 83(9):28077–28123.

773 Nawaf Alampara, Indrajeet Mandal, Pranav Khetarpal,
774 Hargun Singh Grover, Mara Schilling-Wilhelmi,
775 NM Anoop Krishnan, and Kevin Maik Jablonka.
776 Macbench: A multimodal chemistry and materials
777 science benchmark.

778 Nawaf Alampara, Mara Schilling-Wilhelmi, Mar-
779 tiño Ríos-García, Indrajeet Mandal, Pranav
780 Khetarpal, Hargun Singh Grover, NM Krishnan,
781 and Kevin Maik Jablonka. 2024. Probing the
782 limitations of multimodal language models for
783 chemistry and materials research. *arXiv preprint
784 arXiv:2411.16955*.

785 Alon Albalak, Yanai Elazar, Sang Michael Xie, Shayne
786 Longpre, Nathan Lambert, Xinyi Wang, Niklas
787 Muennighoff, Bairu Hou, Liangming Pan, Hae-
788 won Jeong, and 1 others. 2024. A survey on
789 data selection for language models. *arXiv preprint
790 arXiv:2402.16827*.

791 Rawan AlSaad, Alaa Abd-Alrazaq, Sabri Boughorbel,
792 Arfan Ahmed, Max-Antoine Renault, Rafat Damseh,
793 and Javaid Sheikh. 2024. Multimodal large lan-
794 guage models in health care: applications, chal-
795 lenges, and future outlook. *Journal of medical In-
796 ternet research*, 26:e59505.

797 AI Anthropic. 2024. Claude 3.5 sonnet model card ad-
798 dendum. *Claude-3.5 Model Card*, 3.

799 Daman Arora, Himanshu Gaurav Singh, and 1 others.
800 2023. Have llms advanced enough? a challenging
801 problem solving benchmark for large language mod-
802 els. *arXiv preprint arXiv:2305.15074*.

803 Zhangir Azerbayev, Hailey Schoelkopf, Keiran Paster,
804 Marco Dos Santos, Stephen Marcus McAleer, Al-
805 bert Q. Jiang, Jia Deng, Stella Biderman, and Sean
806 Welleck. 2023. **Llemma: An open language model
807 for mathematics.** *ArXiv*, abs/2310.10631.

808 Viraj Bagal, Rishal Aggarwal, P. K. Vinod, and U. Deva
809 Priyakumar. 2021. **Molgpt: Molecular generation
810 using a transformer-decoder model.** *Journal of
811 chemical information and modeling*.

812 Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang,
813 Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou,
814 and Jingren Zhou. 2023. Qwen-vl: A versatile
815 vision-language model for understanding, localiza-
816 tion, text reading, and beyond. *arXiv preprint
817 arXiv:2308.12966*.

818 Tianyi Bai, Hao Liang, Binwang Wan, Yanran Xu, 871
819 Xi Li, Shiyu Li, Ling Yang, Bozhou Li, Yifan Wang, 872
820 Bin Cui, and 1 others. 2024a. A survey of multi- 873
821 modal large language model from a data-centric per- 874
822 spective. *arXiv preprint arXiv:2405.16640*. 875

823 Zechen Bai, Pichao Wang, Tianjun Xiao, Tong He, 876
824 Zongbo Han, Zheng Zhang, and Mike Zheng 877
825 Shou. 2024b. Hallucination of multimodal large 878
826 language models: A survey. *arXiv preprint 879*
827 *arXiv:2404.18930*.

828 Lei Bao, Tianfan Cai, Kathy Koenig, Kai Fang, Jing 880
829 Han, Jing Wang, Qing Liu, Lin Ding, Lili Cui, Ying 881
830 Luo, and 1 others. 2009. Learning and scientific rea- 882
831 soning. *Science*, 323(5914):586–587. 883

832 Kristian G Barman, Sascha Caron, Emily Sullivan, 884
833 Henk W de Regt, Roberto Ruiz de Austri, Mieke 885
834 Boon, Michael Färber, Stefan Fröse, Faegheh Hasibi, 886
835 Andreas Ipp, and 1 others. 2025. Large physics mod- 887
836 els: Towards a collaborative approach with large 888
837 language models and foundation models. *arXiv 889*
838 *preprint arXiv:2501.05382*. 890

839 Khaled Bayouhdh, Raja Knani, Fayçal Hamdaoui, and 891
840 Abdellatif Mtibaa. 2022. A survey on deep mul- 892
841 timodal learning for computer vision: advances, 893
842 trends, applications, and datasets. *The Visual Com- 894*
843 *puter*, 38(8):2939–2970. 895

844 Edward Beeching, Shengyi Costa Huang, Albert Jiang, 896
845 Jia Li, Benjamin Lipkin, Zihan Qina, Kashif Ra- 897
846 sul, Ziju Shen, Roman Soletskyi, and Lewis Tunstall. 898
847 2024. Numinamath 7b tir. [https://huggingface. 899](https://huggingface.co/AI-MO/NuminaMath-7B-TIR)
848 [co/AI-MO/NuminaMath-7B-TIR](https://huggingface.co/AI-MO/NuminaMath-7B-TIR). 900

849 Maciej Besta, Julia Barth, Eric Schreiber, Ales Ku- 901
850 bicek, Afonso Catarino, Robert Gerstenberger, Pi- 902
851 otr Nyczyk, Patrick Iff, Yueling Li, Sam Houliston, 903
852 and 1 others. 2025. Reasoning language models: A 904
853 blueprint. *arXiv preprint arXiv:2501.11223*. 905

854 Abeba Birhane, Atoosa Kasirzadeh, David Leslie, and 906
855 Sandra Wachter. 2023. Science in the age of 907
856 large language models. *Nature Reviews Physics*, 5(5):277–280. 908

858 Erik Cambria, Lorenzo Malandri, Fabio Mercurio, 909
859 Navid Nobani, and Andrea Seveso. 2024. Xai meets 910
860 llms: A survey of the relation between explain- 911
861 able ai and large language models. *arXiv preprint 912*
862 *arXiv:2407.15248*. 913

863 He Cao, Zijing Liu, Xingyu Lu, Yuan Yao, and Yu Li. 914
864 2023a. Instructmol: Multi-modal integration for 915
865 building a versatile and reliable molecular assistant 916
866 in drug discovery. *arXiv preprint arXiv:2311.16208*. 917

867 He Cao, Zijing Liu, Xingyu Lu, Yuan Yao, and Yu Li. 918
868 2023b. **Instructmol: Multi-modal integration for 919**
869 **building a versatile and reliable molecular assistant 920**
870 **in drug discovery**. *ArXiv*, abs/2311.16208. 921

Neeloy Chakraborty, Melkior Ornik, and Katherine 922
Driggs-Campbell. 2024. Hallucination detection in 923
foundation models for decision-making: A flexible 924
definition and review of the state of the art. *arXiv 925*
preprint arXiv:2403.16527. 926

Jinho Chang and Jong-Chul Ye. 2022. **Bidirectional 927**
generation of structure and properties through a sin- 928
gle molecular foundation model. *Nature Communi- 929*
cations, 15. 930

Hila Chefer, Shir Gur, and Lior Wolf. 2021. Generic 931
attention-model explainability for interpreting bi- 932
modal and encoder-decoder transformers. In *Pro- 933*
ceedings of the IEEE/CVF International Conference 934
on Computer Vision, pages 397–406. 935

Qiguang Chen, Libo Qin, Jinhao Liu, Dengyun Peng, 936
Jiannan Guan, Peng Wang, Mengkang Hu, Yuhang 937
Zhou, Te Gao, and Wanxiang Che. 2025. Towards 938
reasoning era: A survey of long chain-of-thought 939
for reasoning large language models. *arXiv preprint 940*
arXiv:2503.09567. 941

Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, 942
Chenfei Yuan, Chi-Min Chan, Heyang Yu, Yaxi 943
Lu, Yi-Hsin Hung, Chen Qian, and 1 others. 2023a. 944
Agentverse: Facilitating multi-agent collaboration 945
and exploring emergent behaviors. In *The Twelfth 946*
International Conference on Learning Representa- 947
tions. 948

Wenhu Chen, Ming Yin, Max Ku, Pan Lu, Yixin Wan, 949
Xueguang Ma, Jianyu Xu, Xinyi Wang, and Tony 950
Xia. 2023b. Theoremqa: A theorem-driven question 951
answering dataset. In *Proceedings of the 2023 Con- 952*
ference on Empirical Methods in Natural Language 953
Processing, pages 7889–7901. 954

Yong Chen, Hongpeng Chen, and Songzhi Su. 2023c. 955
Fine-tuning large language models in education. In 956
2023 13th International Conference on Informa- 957
tion Technology in Medicine and Education (ITME), 958
pages 718–723. IEEE. 959

Zhe Chen, Weiyun Wang, Hao Tian, Shenglong 960
Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, 961
Kongzhi Hu, Jiapeng Luo, Zheng Ma, and 1 oth- 962
ers. 2024a. How far are we to gpt-4v? closing the 963
gap to commercial multimodal models with open- 964
source suites. *Science China Information Sciences*, 965
67(12):220101. 966

Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo 967
Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, 968
Xizhou Zhu, Lewei Lu, and 1 others. 2024b. In- 969
ternvl: Scaling up vision foundation models and 970
aligning for generic visual-linguistic tasks. In *Pro- 971*
ceedings of the IEEE/CVF Conference on Com- 972
puter Vision and Pattern Recognition, pages 24185– 973
24198. 974

Ziru Chen, Shijie Chen, Yuting Ning, Qianheng Zhang, 975
Boshi Wang, Botao Yu, Yifei Li, Zeyi Liao, Chen 976
Wei, Zitong Lu, and 1 others. 2024c. Scienceagent- 977
bench: Toward rigorous assessment of language 978

- 1040 Chen Gao, Xiaochong Lan, Nian Li, Yuan Yuan, Jing- 1097
1041 tao Ding, Zhilun Zhou, Fengli Xu, and Yong Li. 1098
1042 2024a. Large language models empowered agent- 1099
1043 based modeling and simulation: A survey and per- 1100
1044 spectives. *Humanities and Social Sciences Commu- 1101
1045 nications*, 11(1):1–24.
- 1046 Jensen Gao, Bidipta Sarkar, Fei Xia, Ted Xiao, Ji- 1102
1047 ajun Wu, Brian Ichter, Anirudha Majumdar, and 1103
1048 Dorsa Sadigh. 2024b. Physically grounded vision- 1104
1049 language models for robotic manipulation. In *2024 1105
1050 IEEE International Conference on Robotics and Au- 1106
1051 tomation (ICRA)*, pages 12462–12469. IEEE. 1107
1108
- 1052 Steven N Goodman. 2016. Aligning statistical and sci- 1109
1053 entific reasoning. *Science*, 352(6290):1180–1181. 1110
- 1054 Shenghui Guan and Guanyu Wang. 2024. Drug dis- 1111
1055 covery and development in the era of artificial in- 1112
1056 telligence: From machine learning to large lan- 1113
1057 guage models. *Artificial Intelligence Chemistry*, 1114
1058 2(1):100070.
- 1059 Xinyan Guan, Yanjiang Liu, Hongyu Lin, Yaojie Lu, 1115
1060 Ben He, Xianpei Han, and Le Sun. 2024. Mitigating 1116
1061 large language model hallucinations via autonomous 1117
1062 knowledge graph-based retrofitting. In *Proceedings 1118
1063 of the AAAI Conference on Artificial Intelligence*, 1119
1064 volume 38, pages 18126–18134.
- 1065 Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai 1120
1066 Dong, Wentao Zhang, Guanting Chen, Xiao Bi, 1121
1067 Yu Wu, YK Li, and 1 others. 2024a. Deepseek- 1122
1068 coder: When the large language model meets 1123
1069 programming—the rise of code intelligence. *arXiv 1124
1070 preprint arXiv:2401.14196*.
- 1071 Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi 1125
1072 Chang, Shichao Pei, Nitesh V Chawla, Olaf Wiest, 1126
1073 and Xiangliang Zhang. 2024b. Large language 1127
1074 model based multi-agents: A survey of progress and 1128
1075 challenges. *arXiv preprint arXiv:2402.01680*.
- 1076 Taicheng Guo, Bozhao Nan, Zhenwen Liang, Zhichun 1129
1077 Guo, Nitesh Chawla, Olaf Wiest, Xiangliang Zhang, 1130
1078 and 1 others. 2023. What can large language mod- 1131
1079 els do in chemistry? a comprehensive benchmark 1132
1080 on eight tasks. *Advances in Neural Information Pro- 1133
1081 cessing Systems*, 36:59662–59688.
- 1082 Ziyu Guo, Ray Zhang, Hao Chen, Jialin Gao, Dongzhi 1134
1083 Jiang, Jiaze Wang, and Pheng-Ann Heng. 2025a. 1135
1084 **Sciverse: Unveiling the knowledge comprehension 1136
1085 and visual reasoning of llms on multi-modal scien- 1137
1086 tific problems.** *ArXiv*, abs/2503.10627.
- 1087 Ziyu Guo, Renrui Zhang, Hao Chen, Jialin Gao, Hong- 1138
1088 sheng Li, and Pheng-Ann Heng. 2025b. **Sciverse: 1139
1089 Unveiling the knowledge comprehension and visual 1140
1090 reasoning of llms on multi-modal scientific prob- 1141
1091 lems.** *arXiv preprint*.
- 1092 Xixuan Hao, Wei Chen, Yibo Yan, Siru Zhong, Kun 1142
1093 Wang, Qingsong Wen, and Yuxuan Liang. 2024. 1143
1094 Urbanvlp: A multi-granularity vision-language pre- 1144
1095 trained foundation model for urban indicator predic- 1145
1096 tion. *arXiv preprint arXiv:2403.16831*.
- Yunzhuo Hao, Jiawei Gu, Huichen Will Wang, Linjie 1146
Li, Zhengyuan Yang, Lijuan Wang, and Yu Cheng. 1147
2025. Can llms reason in multimodality? emma: 1148
An enhanced multimodal reasoning benchmark. 1149
arXiv preprint arXiv:2501.05444. 1150
1151
- Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding 1152
Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu, 1153
Xu Han, Yujie Huang, Yuxiang Zhang, and 1 oth- 1154
ers. 2024a. Olympiadbench: A challenging bench- 1155
mark for promoting agi with olympiad-level bilin- 1156
gual multimodal scientific problems. *arXiv preprint 1157
arXiv:2402.14008*. 1158
- Xuehai He, Weixi Feng, Kaizhi Zheng, Yujie Lu, Wan- 1159
rong Zhu, Jiachen Li, Yue Fan, Jianfeng Wang, 1160
Linjie Li, Zhengyuan Yang, and 1 others. 2024b. 1161
Mmworld: Towards multi-discipline multi-faceted 1162
world model evaluation in videos. *arXiv preprint 1163
arXiv:2406.08407*. 1164
- Zheqi He, Xinya Wu, Pengfei Zhou, Richeng Xuan, 1165
Guang Liu, Xi Yang, Qiannan Zhu, and Hua Huang. 1166
2024c. Cmmu: A benchmark for chinese multi- 1167
modal multi-type question understanding and rea- 1168
soning. *arXiv preprint arXiv:2401.14011*. 1169
- Dan Hendrycks, Collin Burns, Steven Basart, Andy 1170
Zou, Mantas Mazeika, Dawn Song, and Jacob Stein- 1171
hardt. 2021. Measuring massive multitask language 1172
understanding. *arXiv preprint arXiv:2009.03300*. 1173
- Glen M Hocky. 2024. Connecting molecular properties 1174
with plain language. *Nature Machine Intelligence*, 1175
6(3):249–250. 1176
- Sameera Horawalavithana, Sai Munikoti, Ian Stewart, 1177
and Henry Kvinge. 2023. Scitune: Aligning large 1178
language models with scientific multimodal instruc- 1179
tions. *arXiv preprint arXiv:2307.01139*. 1180
- Bihao Hu, Longwei Zheng, Jiayi Zhu, Lishan Ding, 1181
Yilei Wang, and Xiaoqing Gu. 2024. Teaching plan 1182
generation and evaluation with gpt-4: Unleashing 1183
the potential of llm in instructional design. *IEEE 1184
Transactions on Learning Technologies*. 1185
- Shuaibo Hu and Kui Yu. 2024. Learning robust ratio- 1186
nales for model explainability: A guidance-based 1187
approach. In *Proceedings of the AAAI Conference 1188
on Artificial Intelligence*, volume 38, pages 18243– 1189
18251. 1190
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, 1191
Zhangyin Feng, Haotian Wang, Qianglong Chen, 1192
Weihua Peng, Xiaocheng Feng, Bing Qin, and 1 1193
others. 2023. A survey on hallucination in large 1194
language models: Principles, taxonomy, challenges, 1195
and open questions. *ACM Transactions on Informa- 1196
tion Systems*. 1197
- Yinkui Huang, Tianrun Gao, Jiangjiang Zhang, Xi- 1198
aohong Liu, and Guangyu Wang. 2024a. Adapt- 1199
ing large language models for biomedicine though 1200
retrieval-augmented generation with documents 1201
scoring. In *2024 IEEE International Conference 1202
on Bioinformatics and Biomedicine (ISBI)*. 1203

1153	on <i>Bioinformatics and Biomedicine (BIBM)</i> , pages	John Jumper, Richard Evans, Alexander Pritzel,	1209
1154	5770–5775. IEEE.	Tim Green, Michael Figurnov, Olaf Ronneberger,	1210
1155	Zhen Huang, Zengzhi Wang, Shijie Xia, and Pengfei	Kathryn Tunyasuvunakool, Russ Bates, Augustin	1211
1156	Liu. 2024b. Olympicarena medal ranks: Who is	Žídek, Anna Potapenko, and 1 others. 2021. Highly	1212
1157	the most intelligent ai so far? <i>arXiv preprint</i>	accurate protein structure prediction with alphafold.	1213
1158	<i>arXiv:2406.16772</i> .	<i>nature</i> , 596(7873):583–589.	1214
1159	Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang,	Oren Kraus, Kian Kenyon-Dean, Saber Saberian,	1215
1160	Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun	Maryam Fallah, Peter McLean, Jess Leung, Vasudev	1216
1161	Zhang, Bowen Yu, Keming Lu, and 1 others. 2024.	Sharma, Ayla Khan, Jia Balakrishnan, Safiye Celik,	1217
1162	Qwen2. 5-coder technical report. <i>arXiv preprint</i>	and 1 others. 2024. Masked autoencoders for mi-	1218
1163	<i>arXiv:2409.12186</i> .	croscopy are scalable learners of cellular biology. In	1219
1164	Jiahao Huo, Yibo Yan, Boren Hu, Yutao Yue, and	<i>Proceedings of the IEEE/CVF Conference on Com-</i>	1220
1165	Xuming Hu. 2024. Mmneuron: Discovering	<i>puter Vision and Pattern Recognition</i> , pages 11757–	1221
1166	neuron-level domain-specific interpretation in mul-	11768.	1222
1167	timodal large language model. <i>arXiv preprint</i>	Jakub Lála, Odhran O’Donoghue, Aleksandar Shtedrit-	1223
1168	<i>arXiv:2406.11193</i> .	ski, Sam Cox, Samuel G Rodrigues, and Andrew D	1224
1169	Md Farhan Ishmam, Md Sakib Hossain Shovon,	White. 2023. Paperqa: Retrieval-augmented gener-	1225
1170	Muhammad Firoz Mridha, and Nilanjan Dey. 2024.	ative agent for scientific research. <i>arXiv preprint</i>	1226
1171	From image to language: A critical analysis of	<i>arXiv:2312.07559</i> .	1227
1172	visual question answering (vqa) approaches, chal-	Guillaume Lample, Marie-Anne Lachaux, Thibaut	1228
1173	lenges, and opportunities. <i>Information Fusion</i> , page	Lavril, Xavier Martinet, Amaury Hayat, Gabriel	1229
1174	102270.	Ebner, Aurélien Rodriguez, and Timothée Lacroix.	1230
1175	Aryan Jadon, Avinash Patil, and Shashank Ku-	2022. Hypertree proof search for neural theorem	1231
1176	mar. 2025. Enhancing domain-specific retrieval-	proving . <i>ArXiv</i> , abs/2205.11491.	1232
1177	augmented generation: Synthetic data generation	Antone E Lawson. 2004. The nature and development	1233
1178	and evaluation using reasoning models. <i>arXiv</i>	of scientific reasoning: A synthetic view. <i>Internat-</i>	1234
1179	<i>preprint arXiv:2502.15854</i> .	<i>tional Journal of Science and Mathematics Educa-</i>	1235
1180	Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richard-	<i>tion</i> , 2:307–338.	1236
1181	son, Ahmed El-Kishky, Aiden Low, Alec Helyar,	Martha Lewis and Melanie Mitchell. 2024. Evaluating	1237
1182	Aleksander Madry, Alex Beutel, Alex Carney, and 1	the robustness of analogical reasoning in large lan-	1238
1183	others. 2024. Openai o1 system card. <i>arXiv preprint</i>	guage models. <i>arXiv preprint arXiv:2411.14215</i> .	1239
1184	<i>arXiv:2412.16720</i> .	Aitor Lewkowycz, Anders Andreassen, David Dohan,	1240
1185	Raj Jaiswal, Dhruv Jain, Harsh Parimal Popat, Avinash	Ethan Dyer, Henryk Michalewski, Vinay Venkatesh	1241
1186	Anand, Abhishek Dharmadhikari, Atharva Marathe,	Ramasesh, Ambrose Slone, Cem Anil, Imanol	1242
1187	and Rajiv Ratn Shah. 2024. Improving physics rea-	Schlag, Theo Gutman-Solo, Yuhuai Wu, Behnam	1243
1188	soning in large language models using mixture of re-	Neyshabur, Guy Gur-Ari, and Vedant Misra. 2022.	1244
1189	finement agents. <i>arXiv preprint arXiv:2412.00821</i> .	Solving quantitative reasoning problems with lan-	1245
1190	Xuhui Jiang, Yuxing Tian, Fengrui Hua, Chengjin Xu,	guage models . <i>ArXiv</i> , abs/2206.14858.	1246
1191	Yuanzhuo Wang, and Jian Guo. 2024a. A survey on	Bangzheng Li, Ben Zhou, Fei Wang, Xingyu Fu, Dan	1247
1192	large language model hallucination via a creativity	Roth, and Muhao Chen. 2024a. Deceptive seman-	1248
1193	perspective. <i>arXiv preprint arXiv:2402.06647</i> .	tic shortcuts on reasoning chains: How far can mod-	1249
1194	Zhihuan Jiang, Zhen Yang, Jinhao Chen, Zhengxiao	els go without hallucination? In <i>Proceedings of the</i>	1250
1195	Du, Weihang Wang, Bin Xu, Yuxiao Dong, and Jie	<i>2024 Conference of the North American Chapter of</i>	1251
1196	Tang. 2024b. Visscience: An extensive benchmark	<i>the Association for Computational Linguistics: Hu-</i>	1252
1197	for evaluating k12 educational multi-modal scienti-	<i>man Language Technologies (Volume 1: Long Pa-</i>	1253
1198	fic reasoning. <i>arXiv preprint arXiv:2409.13730</i> .	<i>pers)</i> , pages 7668–7681.	1254
1199	Shiyu Jin, Jinxuan Xu, Yutian Lei, and Liangjun Zhang.	Hang Li, Tianlong Xu, Kaiqi Yang, Yucheng Chu, Yan-	1255
1200	2024a. Reasoning grasping via multimodal large	ling Chen, Yichi Song, Qingsong Wen, and Hui Liu.	1256
1201	language model. <i>arXiv preprint arXiv:2402.06798</i> .	2024b. Ask-before-detection: Identifying and mit-	1257
1202	Zhijing Jin, Yuen Chen, Felix Leeb, Luigi Gresele,	igating conformity bias in llm-powered error detec-	1258
1203	Ojasv Kamal, Zhiheng Lyu, Kevin Blin, Fernando	tor for math word problem solutions. <i>arXiv preprint</i>	1259
1204	Gonzalez Aduato, Max Kleiman-Weiner, Mrinmaya	<i>arXiv:2412.16838</i> .	1260
1205	Sachan, and 1 others. 2024b. Cladder: A benchmark	Hang Li, Tianlong Xu, Chaoli Zhang, Eason Chen,	1261
1206	to assess causal reasoning capabilities of language	Jing Liang, Xing Fan, Haoyang Li, Jiliang Tang,	1262
1207	models. <i>Advances in Neural Information Process-</i>	and Qingsong Wen. 2024c. Bringing generative ai	1263
1208	<i>ing Systems</i> , 36.	to adaptive learning in education. <i>arXiv preprint</i>	1264
		<i>arXiv:2402.14601</i> .	1265

1266	Haonan Li, Yixuan Zhang, Fajri Koto, Yifei Yang, Hai Zhao, Yeyun Gong, Nan Duan, and Timothy Baldwin. 2024d. CMMLU: Measuring massive multitask language understanding in Chinese . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 11260–11285, Bangkok, Thailand. Association for Computational Linguistics.	1319
1267		1320
1268		1321
1269		1322
1270		1323
1271		1324
1272		
1273	Jian Li, Weiheng Lu, Hao Fei, Meng Luo, Ming Dai, Min Xia, Yizhang Jin, Zhenye Gan, Ding Qi, Chaoyou Fu, and 1 others. 2024e. A survey on benchmarks of multimodal large language models. <i>arXiv preprint arXiv:2408.08632</i> .	1325
1274		1326
1275		1327
1276		1328
1277		1329
1278		1330
1279		
1280	Junyi Li, Jie Chen, Ruiyang Ren, Xiaoxue Cheng, Wayne Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. 2024f. The dawn after the dark: An empirical study on factuality hallucination in large language models. <i>arXiv preprint arXiv:2401.03205</i> .	1331
1281		1332
1282		1333
1283		1334
1284		1335
1285		1336
1286		1337
1287		1338
1288		1339
1289		1340
1290		1341
1291		
1292		
1293	Lei Li, Yuqi Wang, Runxin Xu, Peiyi Wang, Xiachong Feng, Lingpeng Kong, and Qi Liu. 2024h. Multimodal arxiv: A dataset for improving scientific comprehension of large vision-language models. <i>arXiv preprint arXiv:2403.00231</i> .	1342
1294		1343
1295		1344
1296		1345
1297		1346
1298		
1299		
1300		
1301		
1302		
1303		
1304		
1305		
1306		
1307		
1308		
1309		
1310		
1311		
1312		
1313		
1314		
1315		
1316		
1317		
1318		
1319	Zekun Li, Xianjun Yang, Kyuri Choi, Wanrong Zhu, Ryan Hsieh, HyeonJung Kim, Jin Hyuk Lim, Sungyoun Ji, Byungju Lee, Xifeng Yan, and 1 others. 2024m. Mmsci: A multimodal multi-discipline dataset for phd-level scientific comprehension. In <i>AI for Accelerated Materials Design-Vienna 2024</i> .	1320
1320		1321
1321		1322
1322		1323
1323		1324
1324		
1325		
1326		
1327		
1328		
1329		
1330		
1331		
1332		
1333		
1334		
1335		
1336		
1337		
1338		
1339		
1340		
1341		
1342		
1343		
1344		
1345		
1346		
1347		
1348		
1349		
1350		
1351		
1352		
1353		
1354		
1355		
1356		
1357		
1358		
1359		
1360		
1361		
1362		
1363		
1364		
1365		
1366		
1367		
1368		
1369		
1370		
1371		
1372		
1373		

1374	mathematics reasoning of large multimodal models.	Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jian-	1430
1375	<i>arXiv preprint arXiv:2409.02834</i> .	Guang Lou, Chongyang Tao, Xiubo Geng, Qingwei	1431
1376	Xiang Liu, Penglei Sun, Shuyan Chen, Longhan	Lin, Shifeng Chen, and Dongmei Zhang. 2023. Wiz-	1432
1377	Zhang, Peijie Dong, Huajie You, Yongqi Zhang,	ardmath: Empowering mathematical reasoning for	1433
1378	Chang Yan, Xiaowen Chu, and Tong yi Zhang.	large language models via reinforced evol-instruct.	1434
1379	2025c. Perovskite-llm: Knowledge-enhanced large	<i>ArXiv</i> , abs/2308.09583.	1435
1380	language models for perovskite solar cell research.	Renqian Luo, Liai Sun, Yingce Xia, Tao Qin, Sheng	1436
1381	<i>Preprint</i> , arXiv:2502.12669.	Zhang, Hoifung Poon, and Tie-Yan Liu. 2022.	1437
1382	Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li,	Biogpt: Generative pre-trained transformer for	1438
1383	Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi	biomedical text generation and mining. <i>Briefings in</i>	1439
1384	Wang, Conghui He, Ziwei Liu, and 1 others. 2025d.	<i>bioinformatics</i> .	1440
1385	Mmbench: Is your multi-modal model an all-around	Rachel K. Luu and M. Buehler. 2023. Bioinspiredllm:	1441
1386	player? In <i>European conference on computer vision</i> ,	Conversational large language model for the me-	1442
1387	pages 216–233. Springer.	chanics of biological and bioinspired materials. <i>Ad-</i>	1443
1388	Micha Livne, Zulfat Miftahutdinov, E. Tutubalina,	<i>vanced Science</i> , 11.	1444
1389	Maksim Kuznetsov, Daniil Polykovskiy, Annika	Liuzhenghao Lv, Zongying Lin, Hao Li, Yuyang Liu,	1445
1390	Brundyn, Aastha Jhunjhunwala, Anthony B Costa,	Jiaxi Cui, Calvin Yu-Chian Chen, Li Yuan, and	1446
1391	Alex Aliper, and Alex Zhavoronkov. 2023. nach0:	Yonghong Tian. 2024. Prollama: A protein lan-	1447
1392	multimodal natural and chemical languages founda-	guage model for multi-task protein language pro-	1448
1393	tion model. <i>Chemical Science</i> , 15:8380 – 8389.	cessing. <i>ArXiv</i> .	1449
1394	Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foer-	Andres M. Bran, Sam Cox, Oliver Schilter, Carlo Bal-	1450
1395	ster, Jeff Clune, and David Ha. 2024a. The ai scient-	dassari, Andrew D White, and Philippe Schwaller.	1451
1396	ist: Towards fully automated open-ended scientific	2024. Augmenting large language models with	1452
1397	discovery. <i>arXiv preprint arXiv:2408.06292</i> .	chemistry tools. <i>Nature Machine Intelligence</i> , pages	1453
1398	Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chun-	1–11.	1454
1399	yuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-	Xueguang Ma, Qian Liu, Dongfu Jiang, Ge Zhang,	1455
1400	wei Chang, Michel Galley, and Jianfeng Gao. 2023.	Zejun Ma, and Wenhu Chen. 2025. General-	1456
1401	Mathvista: Evaluating math reasoning in visual con-	reasoner: Advancing llm reasoning across all do-	1457
1402	texts with gpt-4v, bard, and other large multimodal	mains. <i>Preprint</i> , arXiv:2505.14652.	1458
1403	models. <i>arXiv e-prints</i> , pages arXiv–2310.	Gengchen Mai, Weiming Huang, Jin Sun, Suhang	1459
1404	Pan Lu, Ran Gong, Shibiao Jiang, Liang Qiu, Siyuan	Song, Deepak Mishra, Ninghao Liu, Song Gao,	1460
1405	Huang, Xiaodan Liang, and Song-Chun Zhu. 2021a.	Tianming Liu, Gao Cong, Yingjie Hu, and 1 others.	1461
1406	Inter-gps: Interpretable geometry problem solving	2024. On the opportunities and challenges of founda-	1462
1407	with formal language and symbolic reasoning. <i>arXiv</i>	tion models for geoai (vision paper). <i>ACM Trans-</i>	1463
1408	<i>preprint arXiv:2105.04165</i> .	<i>actions on Spatial Algorithms and Systems</i> .	1464
1409	Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-	Roger G Melko and Juan Carrasquilla. 2024. Language	1465
1410	Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter	models for quantum simulation. <i>Nature Computa-</i>	1466
1411	Clark, and Ashwin Kalyan. 2022a. Learn to explain:	<i>tional Science</i> , 4(1):11–18.	1467
1412	Multimodal reasoning via thought chains for science	Bonan Min, Hayley Ross, Elior Sulem, Amir	1468
1413	question answering. <i>Advances in Neural Informa-</i>	Pouran Ben Veyseh, Thien Huu Nguyen, Oscar	1469
1414	<i>tion Processing Systems</i> , 35:2507–2521.	Sainz, Eneko Agirre, Ilana Heintz, and Dan Roth.	1470
1415	Pan Lu, Liang Qiu, Jiaqi Chen, Tony Xia, Yizhou Zhao,	2023. Recent advances in natural language process-	1471
1416	Wei Zhang, Zhou Yu, Xiaodan Liang, and Song-	ing via large pre-trained language models: A survey.	1472
1417	Chun Zhu. 2021b. Iconqa: A new benchmark for	<i>ACM Computing Surveys</i> , 56(2):1–40.	1473
1418	abstract diagram understanding and visual language	Santiago Miret and Nandan M Krishnan. 2024. Are	1474
1419	reasoning. <i>arXiv preprint arXiv:2110.13214</i> .	llms ready for real-world materials discovery?	1475
1420	Pan Lu, Liang Qiu, Wenhao Yu, Sean Welleck, and	<i>arXiv preprint arXiv:2402.05200</i> .	1476
1421	Kai-Wei Chang. 2022b. A survey of deep learn-	Adrian Mirza, Nawaf Alampara, Sreekanth Kunchapu,	1477
1422	ing for mathematical reasoning. <i>arXiv preprint</i>	Martiño Ríos-García, Benedict Emoekabu, Aswanth	1478
1423	<i>arXiv:2212.10535</i> .	Krishnan, Tanya Gupta, Mara Schilling-Wilhelmi,	1479
1424	Zimu Lu, Aojun Zhou, Ke Wang, Houxing Ren,	Macjonathan Okereke, Anagha Aneesh, and 1 oth-	1480
1425	Weikang Shi, Juntong Pan, Mingjie Zhan, and	ers. 2024. Are large language models superhuman	1481
1426	Hongsheng Li. 2024b. Mathcoder2: Better math	chemists? <i>arXiv preprint arXiv:2404.01475</i> .	1482
1427	reasoning from continued pretraining on model-		
1428	translated mathematical code. <i>arXiv preprint</i>		
1429	<i>arXiv:2410.08196</i> .		

1483	Iman Mirzadeh, Keivan Alizadeh, Hooman Shahrokhi, Oncel Tuzel, Samy Bengio, and Mehrdad Farajtabar. 2024. Gsm-symbolic: Understanding the limitations of mathematical reasoning in large language models. <i>arXiv preprint arXiv:2410.05229</i> .	1534
1484		1535
1485		1536
1486		1537
1487		1538
1488	Mistral AI Team. 2024. Mathstral .	1539
1489	MoonshotAI. 2024. k0-math official platform .	
1490	Albert Morera. 2024. Foundation models in shaping the future of ecology. <i>Ecological Informatics</i> , page 102545.	
1491		
1492		
1493	Meredith Ringel Morris. 2023. Scientists’ perspectives on the potential for generative ai in their fields. <i>arXiv preprint arXiv:2304.01420</i> .	
1494		
1495		
1496	Spyridon Mouselinos, Henryk Michalewski, and Mateusz Malinowski. 2024. Beyond lines and circles: Unveiling the geometric reasoning gap in large language models. <i>arXiv preprint arXiv:2402.03877</i> .	
1497		
1498		
1499		
1500	Niklas Muennighoff, Alexander Rush, Boaz Barak, Teven Le Scao, Nouamane Tazi, Aleksandra Piktus, Sampo Pyysalo, Thomas Wolf, and Colin A Raffel. 2023. Scaling data-constrained language models. <i>Advances in Neural Information Processing Systems</i> , 36:50358–50376.	
1501		
1502		
1503		
1504		
1505		
1506	Nikahat Mulla and Prachi Gharpure. 2023. Automatic question generation: a review of methodologies, datasets, evaluation metrics, and applications. <i>Progress in Artificial Intelligence</i> , 12(1):1–32.	
1507		
1508		
1509		
1510	Alhassan Mumuni and Fuseini Mumuni. 2025. Large language models for artificial general intelligence (agi): A survey of foundational principles and approaches. <i>arXiv preprint arXiv:2501.03151</i> .	
1511		
1512		
1513		
1514	Siddharth M. Narayanan, James D. Braza, Ryan-Rhys Griffiths, Albert Bou, Geemi Wellawatte, Mayk Caldas Ramos, Ludovico Mitchener, Samuel G. Rodrigues, and Andrew D. White. 2025. Training a scientific reasoning model for chemistry . <i>Preprint</i> , arXiv:2506.17238.	
1515		
1516		
1517		
1518		
1519		
1520	Tung Nguyen, Johannes Brandstetter, Ashish Kapoor, Jayesh K Gupta, and Aditya Grover. 2023. Climax: A foundation model for weather and climate. <i>arXiv preprint arXiv:2301.10343</i> .	
1521		
1522		
1523		
1524	Bolin Ni, Jingcheng Hu, Yixuan Wei, Houwen Peng, Zheng Zhang, Gaofeng Meng, and Han Hu. 2024. Xwin-llm: Strong and scalable alignment practice for llms . <i>ArXiv</i> , abs/2405.20335.	
1525		
1526		
1527		
1528	Mengjia Niu, Hao Li, Jie Shi, Hamed Haddadi, and Fan Mo. 2024. Mitigating hallucinations in large language models via self-refinement-enhanced knowledge retrieval. <i>arXiv preprint arXiv:2405.06545</i> .	
1529		
1530		
1531		
1532	Skywork o1 Team. 2024. Skywork-o1 open series . https://huggingface.co/Skywork .	
1533		
	David Oniani, Jordan Hilsman, Chengxi Zang, Junmei Wang, Lianjin Cai, Jan Zawala, and Yanshan Wang. 2024. Emerging opportunities of using large language language models for translation between drug molecules and indications. <i>arXiv preprint arXiv:2402.09588</i> .	1540
		1541
		1542
		1543
	Peter Onu, Anup Pradhan, and Charles Mbohwa. 2024. Potential to use metaverse for future teaching and learning. <i>Education and Information Technologies</i> , 29(7):8893–8924.	1544
		1545
		1546
	Liangming Pan, Michael Saxon, Wenda Xu, Deepak Nathani, Xinyi Wang, and William Yang Wang. 2023. Automatically correcting large language models: Surveying the landscape of diverse self-correction strategies. <i>arXiv preprint arXiv:2308.03188</i> .	1547
		1548
		1549
	Shirui Pan, Linhao Luo, Yufei Wang, Chen Chen, Jiapu Wang, and Xindong Wu. 2024. Unifying large language models and knowledge graphs: A roadmap. <i>IEEE Transactions on Knowledge and Data Engineering</i> .	1550
		1551
		1552
		1553
		1554
	Priyaranjan Pattanayak, Hitesh Laxmichand Patel, Bhargava Kumar, Amit Agarwal, Ishan Banerjee, Srikant Panda, and Tejaswini Kumar. 2024. Survey of large multimodal model datasets, application categories and taxonomy. <i>arXiv preprint arXiv:2412.17759</i> .	1555
		1556
		1557
		1558
		1559
	Qizhi Pei, Lijun Wu, Kaiyuan Gao, Jinhua Zhu, Yue Wang, Zun Wang, Tao Qin, and Rui Yan. 2024. Leveraging biomolecule and natural language through multi-modal learning: A survey. <i>arXiv preprint arXiv:2403.01528</i> .	1560
		1561
		1562
		1563
		1564
	Alexander R Pelletier, Joseph Ramirez, Irsyad Adam, Simha Sankar, Yu Yan, Ding Wang, Dylan Steinecke, Wei Wang, and Peipei Ping. 2024. Explainable biomedical hypothesis generation via retrieval augmented generation enabled large language models. <i>arXiv preprint arXiv:2407.12888</i> .	1565
		1566
		1567
		1568
		1569
		1570
	Stanislas Polu and Ilya Sutskever. 2020. Generative language modeling for automated theorem proving . <i>ArXiv</i> , abs/2009.03393.	1571
		1572
		1573
	Vignesh Prabhakar, Md Amirul Islam, Adam Atanas, Yao-Ting Wang, Joah Han, Aastha Jhunjhunwala, Rucha Apte, Robert Clark, Kang Xu, Zihan Wang, and Kai Liu. 2025. Omniscience: A domain-specialized llm for scientific reasoning and discovery . <i>Preprint</i> , arXiv:2503.17604.	1574
		1575
		1576
		1577
		1578
		1579
	Michael H Prince, Henry Chan, Aikaterini Vriza, Tao Zhou, Varuni K Sastry, Yanqi Luo, Matthew T Dearing, Ross J Harder, Rama K Vasudevan, and Mathew J Cherukara. 2024. Opportunities for retrieval and tool augmented large language models in scientific facilities. <i>npj Computational Materials</i> , 10(1):251.	1580
		1581
		1582
		1583
		1584
		1585
		1586
	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. <i>arXiv preprint arXiv:2311.05965</i> .	1587
		1588
		1589
		1590

1591 Biqing Qi, Kaiyan Zhang, Kai Tian, Haoxiang
1592 Li, Zhang-Ren Chen, Sihang Zeng, Ermo Hua,
1593 Hu Jinfang, and Bowen Zhou. 2024. Large lan-
1594 guage models as biomedical hypothesis genera-
1595 tors: a comprehensive evaluation. *arXiv preprint*
1596 *arXiv:2407.08940*.

1597 Runqi Qiao, Qiuna Tan, Guanting Dong, Minhui Wu,
1598 Chong Sun, Xiaoshuai Song, Zhuoma GongQue,
1599 Shanglin Lei, Zhe Wei, Miaoxuan Zhang, and 1 oth-
1600 ers. 2024. We-math: Does your large multimodal
1601 model achieve human-like mathematical reasoning?
1602 *arXiv preprint arXiv:2407.01284*.

1603 Qwen. 2024a. **Qvq: To see the world with wisdom**.

1604 Qwen. 2024b. Qwq: Reflect deeply on the bound-
1605 aries of the unknown, november 2024. URL
1606 <https://qwenlm.github.io/blog/qwq-32b-preview>.

1607 Anil Rahate, Rahee Walambe, Sheela Ramanna, and
1608 Ketan Kotecha. 2022. Multimodal co-learning:
1609 Challenges, applications with datasets, recent ad-
1610 vances and future directions. *Information Fusion*,
1611 81:203–239.

1612 Md Abdur Rahman, Lamyaa Alqahtani, Amna Alboq,
1613 and Alaa Ainousah. 2024. A survey on security and
1614 privacy of large multimodal deep learning models:
1615 Teaching and learning perspective. In *2024 21st*
1616 *Learning and Technology Conference (L&T)*, pages
1617 13–18. IEEE.

1618 Mayk Caldas Ramos, Christopher J Collison, and An-
1619 drew D White. 2025. A review of large language
1620 models and autonomous agents in chemistry. *Chem-*
1621 *ical Science*.

1622 Chandan K Reddy and Parshin Shojaee. 2024. Towards
1623 scientific discovery with generative ai: Progress,
1624 opportunities, and challenges. *arXiv preprint*
1625 *arXiv:2412.11427*.

1626 Zhongwei Ren, Zhicheng Huang, Yunchao Wei, Yao
1627 Zhao, Dongmei Fu, Jiashi Feng, and Xiaojie Jin.
1628 2024. Pixellm: Pixel reasoning with large multi-
1629 modal model. In *Proceedings of the IEEE/CVF Con-*
1630 *ference on Computer Vision and Pattern Recogni-*
1631 *tion*, pages 26374–26383.

1632 Jonathan Roberts, Timo Lüddecke, Rehan Sheikh, Kai
1633 Han, and Samuel Albanie. 2024. Charting new ter-
1634 ritories: Exploring the geographic and geospatial
1635 capabilities of multimodal llms. In *Proceedings of the*
1636 *IEEE/CVF Conference on Computer Vision and Pat-*
1637 *tern Recognition*, pages 554–563.

1638 Nikolaos Rodis, Christos Sardanios, Panagiotis
1639 Radoglou-Grammatikis, Panagiotis Sarigiannidis,
1640 Iraklis Varlamis, and Georgios Th Papadopoulos.
1641 2024. Multimodal explainable artificial intelli-
1642 gence: A comprehensive review of methodological
1643 advances and future research directions. *IEEE*
1644 *Access*.

Keisuke Sato. 2024. Exploring the educational land-
scape of ai: Large language models’ approaches to
explaining conservation of momentum in physics.
arXiv preprint arXiv:2407.05308.

Zhihong Shao, Yuxiang Luo, Chengda Lu, ZZ Ren,
Jiewen Hu, Tian Ye, Zhibin Gou, Shirong Ma, and
Xiaokang Zhang. 2025. Deepseekmath-v2: To-
wards self-verifiable mathematical reasoning. *arXiv*
preprint arXiv:2511.22570.

Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu,
Jun-Mei Song, Mingchuan Zhang, Y. K. Li, Yu Wu,
and Daya Guo. 2024. **Deepseekmath: Pushing the**
limits of mathematical reasoning in open language
models. *ArXiv*, abs/2402.03300.

Freda Shi, Xinyun Chen, Kanishka Misra, Nathan
Scales, David Dohan, Ed H Chi, Nathanael Schärli,
and Denny Zhou. 2023. Large language models can
be easily distracted by irrelevant context. In *Inter-*
national Conference on Machine Learning, pages
31210–31227. PMLR.

Haizhou Shi, Zihao Xu, Hengyi Wang, Weiyi Qin,
Wenyuan Wang, Yibin Wang, Zifeng Wang, Sayna
Ebrahimi, and Hao Wang. 2024a. Continual learn-
ing of large language models: A comprehensive sur-
vey. *arXiv preprint arXiv:2404.16789*.

Wenhao Shi, Zhiqiang Hu, Yi Bin, Junhua Liu, Yang
Yang, See-Kiong Ng, Li Bing, and Roy Ka wei
Lee. 2024b. **Math-llava: Bootstrapping mathemati-**
cal reasoning for multimodal large language mod-
els. In *Conference on Empirical Methods in Natural*
Language Processing.

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

Parshin Shojaee, Ngoc-Hieu Nguyen, Kazem Meidani,
Amir Barati Farimani, Khoa D Doan, and Chan-
dan K Reddy. 2025. Llm-srbench: A new bench-
mark for scientific equation discovery with large lan-
guage models. *arXiv preprint arXiv:2504.10415*.

Ali Shtarbanov, Mengjia Zhu, Nicholas Colonnese,
and Amirhossein Hajiagha Memar. 2023. Sleeveio:
modular and reconfigurable platform for multimodal
wearable haptic feedback interactions. In *Proceed-*
ings of the 36th annual ACM symposium on user in-
terface software and technology, pages 1–15.

Michael D Skarlinski, Sam Cox, Jon M Laurent,
James D Braza, Michaela Hinks, Michael J Ham-
merling, Manvitha Ponnampati, Samuel G Rodrigues,
and Andrew D White. 2024. Language agents
achieve superhuman synthesis of scientific knowl-
edge. *arXiv preprint arXiv:2409.13740*.

Eduardo Soares, Victor Shirasuna, Emilio Vital Brazil,
Renato Cerqueira, Dmitry Zubarev, and Kristin
Schmidt. 2024. **A large encoder-decoder family of**

1701 foundation models for chemical language. *ArXiv*,
1702 abs/2407.20267.

1703 Zhangde Song, Jieyu Lu, Yuanqi Du, Botao Yu,
1704 Thomas M Pruyun, Yue Huang, Kehan Guo, Xiuzhe
1705 Luo, Yuanhao Qu, Yi Qu, and 1 others. 2025. Evalu-
1706 ating large language models in scientific discovery.
1707 *arXiv preprint arXiv:2512.15567*.

1708 Rickard Stureborg, Dimitris Alikaniotis, and Yoshi
1709 Suhara. 2024. Large language models are inconsis-
1710 tent and biased evaluators. *ArXiv*, abs/2405.01724.

1711 Jiankai Sun, Chuanyang Zheng, Enze Xie, Zhengying
1712 Liu, Ruihang Chu, Jianing Qiu, Jiaqi Xu, Mingyu
1713 Ding, Hongyang Li, Mengzhe Geng, and 1 others.
1714 2023. A survey of reasoning with foundation mod-
1715 els. *arXiv preprint arXiv:2312.11562*.

1716 Kai Sun, Yushi Bai, Ji Qi, Lei Hou, and Juanzi Li.
1717 2024a. Mm-math: Advancing multimodal math
1718 evaluation with process evaluation and fine-grained
1719 classification. In *Findings of the Association for*
1720 *Computational Linguistics: EMNLP 2024*, pages
1721 1358–1375.

1722 Liangtai Sun, Yang Han, Zihan Zhao, Da Ma, Zhennan
1723 Shen, Baocai Chen, Lu Chen, and Kai Yu. 2024b.
1724 Scieval: A multi-level large language model eval-
1725 uation benchmark for scientific research. In *Pro-*
1726 *ceedings of the AAAI Conference on Artificial Intel-*
1727 *ligence*, volume 38, pages 19053–19061.

1728 Ash Sze and Soha Hassoun. 2024. Evaluation of
1729 search-enabled pre-trained large language models
1730 on retrieval tasks for the pubchem database. *bioRxiv*,
1731 pages 2024–08.

1732 TALEducation. 2023. **Mathgpt official platform**.

1733 Jiabin Tang, Lianghao Xia, Zhonghang Li, and Chao
1734 Huang. 2025. Ai-researcher: Autonomous scientific
1735 innovation. *arXiv preprint arXiv:2505.18705*.

1736 Xin Tang, Jiawei Zhang, Yichun He, Xinhe Zhang,
1737 Zuwan Lin, Sebastian Partarrieu, Emma Bou Hanna,
1738 Zhaolin Ren, Hao Shen, Yuhong Yang, and 1 oth-
1739 ers. 2023. Explainable multi-task learning for multi-
1740 modality biological data analysis. *Nature communi-*
1741 *cations*, 14(1):2546.

1742 Ross Taylor, Marcin Kardas, Guillem Cucurull,
1743 Thomas Scialom, Anthony Hartshorn, Elvis Saravia,
1744 Andrew Poulton, Viktor Kerkez, and Robert Stojnic.
1745 2022. Galactica: A large language model for sci-
1746 ence. *arXiv preprint arXiv:2211.09085*.

1747 Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan
1748 Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer,
1749 Damien Vincent, Zhufeng Pan, Shibo Wang, and 1
1750 others. 2024a. Gemini 1.5: Unlocking multimodal
1751 understanding across millions of tokens of context.
1752 *arXiv preprint arXiv:2403.05530*.

Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan
Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer,
Damien Vincent, Zhufeng Pan, Shibo Wang, and 1
others. 2024b. Gemini 1.5: Unlocking multimodal
understanding across millions of tokens of context.
arXiv preprint arXiv:2403.05530.

Omkar Thawakar, Dinura Dissanayake, Ketan More,
Ritesh Thawkar, Ahmed Heakl, Noor Ahsan, Yuhao
Li, Mohammed Zumri, Jean Lahoud, Rao Muham-
mad Anwer, Hisham Cholakkal, Ivan Laptev,
Mubarak Shah, Fahad Shahbaz Khan, and Salman
Khan. 2025. **Llamav-o1: Rethinking step-by-**
step visual reasoning in llms. *arXiv preprint*
arXiv:2501.06186.

Shagun Uppal, Sarthak Bhagat, Devamanyu Hazarika,
Navonil Majumder, Soujanya Poria, Roger Zimmer-
mann, and Amir Zadeh. 2022. Multimodal research
in vision and language: A review of current and
emerging trends. *Information Fusion*, 77:149–171.

J Vaishnavi and V Narmatha. 2024. Video captioning—
a survey. *Multimedia Tools and Applications*, pages
1–32.

Benyou Wang, Qianqian Xie, Jiahuan Pei, Zhihong
Chen, Prayag Tiwari, Zhao Li, and Jie Fu. 2023a.
Pre-trained language models in biomedical domain:
A systematic survey. *ACM Computing Surveys*,
56(3):1–52.

Cunxiang Wang, Xiaoze Liu, Yuanhao Yue, Xiangru
Tang, Tianhang Zhang, Cheng Jiayang, Yunzhi Yao,
Wenyang Gao, Xuming Hu, Zehan Qi, and 1 others.
2023b. Survey on factuality in large language mod-
els: Knowledge, retrieval and domain-specificity.
arXiv preprint arXiv:2310.07521.

Feng Wang. 2024. Lighthouse: A survey of agi hallu-
cination. *arXiv preprint arXiv:2401.06792*.

Hongyu Wang, Jiayu Xu, Senwei Xie, Ruiping Wang,
Jialin Li, Zhaojie Xie, Bin Zhang, Chuyan Xiong,
and Xilin Chen. 2024a. M4u: Evaluating multi-
lingual understanding and reasoning for large mul-
timodal models. *arXiv preprint arXiv:2405.15638*.

Ke Wang, Juntong Pan, Linda Wei, Aojun Zhou,
Weikang Shi, Zimu Lu, Han Xiao, Yunqiao Yang,
Houxing Ren, Mingjie Zhan, and 1 others. 2025a.
Mathcoder-vl: Bridging vision and code for en-
hanced multimodal mathematical reasoning. *arXiv*
preprint arXiv:2505.10557.

Ke Wang, Houxing Ren, Aojun Zhou, Zimu Lu,
Sichun Luo, Weikang Shi, Renrui Zhang, Linqi
Song, Mingjie Zhan, and Hongsheng Li. 2023c.
Mathcoder: Seamless code integration in llms
for enhanced mathematical reasoning. *ArXiv*,
abs/2310.03731.

Lei Wang, Yi Hu, Jiabang He, Xing Xu, Ning Liu, Hui
Liu, and Heng Tao Shen. 2024b. T-sciq: Teaching
multimodal chain-of-thought reasoning via large lan-
guage model signals for science question answering.

1809 In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 19162–19170. 1865

1810 2024a. Usable xai: 10 strategies towards exploit- 1866

1811 Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao 1867

1812 Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, 1868

1813 Xu Chen, Yankai Lin, and 1 others. 2024c. A survey 1869

1814 on large language model based autonomous agents. 1870

1815 *Frontiers of Computer Science*, 18(6):186345. 1871

1816 Peiyi Wang, Lei Li, Liang Chen, Zefan Cai, Dawei Zhu, 1872

1817 Binghuai Lin, Yunbo Cao, Lingpeng Kong, Qi Liu, 1873

1818 Tianyu Liu, and Zhifang Sui. 2024d. **Large lan-** 1874

1819 **guage models are not fair evaluators.** In *Proceed-* 1875

1820 *ings of the 62nd Annual Meeting of the Association* 1876

1821 *for Computational Linguistics (Volume 1: Long Pa-*

1822 *pers)*, pages 9440–9450, Bangkok, Thailand. Asso- 1877

1823 ciation for Computational Linguistics. 1878

1824 Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhi- 1879

1825 hao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin 1880

1826 Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei 1881

1827 Du, Xuancheng Ren, Rui Men, Dayiheng Liu, 1882

1828 Chang Zhou, Jingren Zhou, and Junyang Lin. 2024e. 1883

1829 Qwen2-vl: Enhancing vision-language model’s per- 1884

1830 ception of the world at any resolution. *arXiv* 1885

1831 *preprint arXiv:2409.12191.* 1886

1832 Xiaoxuan Wang, Ziniu Hu, Pan Lu, Yanqiao Zhu, Jieyu 1887

1833 Zhang, Satyen Subramaniam, Arjun R Loomba, 1888

1834 Shichang Zhang, Yizhou Sun, and Wei Wang. 1889

1835 2024f. Scibench: Evaluating college-level scientific 1890

1836 problem-solving abilities of large language models. 1891

1837 *arXiv preprint arXiv:2307.10635.* 1892

1838 Yiqi Wang, Wentao Chen, Xiaotian Han, Xudong Lin, 1893

1839 Haiteng Zhao, Yongfei Liu, Bohan Zhai, Jianbo 1894

1840 Yuan, Quanzeng You, and Hongxia Yang. 2024g. 1895

1841 Exploring the reasoning abilities of multimodal 1896

1842 large language models (mllms): A comprehensive 1897

1843 survey on emerging trends in multimodal reasoning. 1898

1844 *arXiv preprint arXiv:2401.06805.* 1899

1845 Yizhou Wang, Chen Tang, Han Deng, Jiabei Xiao, Ji- 1899

1846 aqi Liu, Jianyu Wu, Jun Yao, Pengze Li, Encheng 1900

1847 Su, Lintao Wang, and 1 others. 2025b. Scireasoner: 1901

1848 Laying the scientific reasoning ground across disci- 1902

1849 plines. *arXiv preprint arXiv:2509.21320.* 1903

1850 Taylor Webb, Keith J Holyoak, and Hongjing Lu. 2023. 1904

1851 Emergent analogical reasoning in large language 1905

1852 models. *Nature Human Behaviour*, 7(9):1526– 1906

1853 1541. 1907

1854 Jingxuan Wei, Cheng Tan, Zhangyang Gao, Linzhuang 1908

1855 Sun, Siyuan Li, Bihui Yu, Ruifeng Guo, and Stan Z 1909

1856 Li. 2024. Enhancing human-like multimodal reason- 1910

1857 ing: a new challenging dataset and comprehensive 1911

1858 framework. *Neural Computing and Applications*, 1912

1859 36(33):20849–20861. 1913

1860 Andrew D White. 2023. The future of chemistry is lan- 1914

1861 guage. *Nature Reviews Chemistry*, 7(7):457–458. 1915

1862 Xuansheng Wu, Haiyan Zhao, Yaochen Zhu, Yucheng 1916

1863 Shi, Fan Yang, Tianming Liu, Xiaoming Zhai, Wen- 1917

1864 lin Yao, Jundong Li, Mengnan Du, and 1 others. 1918

2024b. **Molmetalm: a physicochemical** 1868

knowledge-guided molecular meta language model. 1869

ArXiv, abs/2411.15500. 1870

2024c. **Internlm2.5-stepprover: Advancing automated theo-** 1872

rem proving via expert iteration on large-scale lean 1873

problems. *ArXiv*, abs/2410.15700. 1874

2023. **Mole-bert: Rethinking pre-training graph neu-** 1875

ral networks for molecules. In *International Confer-* 1876

ence on Learning Representations. 1877

Wenyi Xiao, Ziwei Huang, Leilei Gan, Wangui He, 1882

Haoyuan Li, Zhelun Yu, Hao Jiang, Fei Wu, and Lin- 1883

chao Zhu. 2024. Detecting and mitigating hallucina- 1884

tion in large vision language models via fine-grained 1885

ai feedback. *arXiv preprint arXiv:2404.14233.* 1886

Tong Xie, Yuwei Wan, Wei Huang, Zhenyu Yin, 1887

Yixuan Liu, Shaozhou Wang, Qingyuan Linghu, 1888

Chunyu Kit, Clara Grazian, Wenjie Zhang, Imran 1889

Razzak, and Bram Hoex. 2023. **Darwin series: Do-** 1890

main specific large language models for natural sci- 1891

ence. *ArXiv*, abs/2308.13565. 1892

Tong Xie, Yuwei Wan, Yixuan Liu, Yuchen Zeng, Wen- 1893

jie Zhang, Chunyu Kit, Dongzhan Zhou, and Bram 1894

Hoex. 2024. **Darwin 1.5: Large language models as** 1895

materials science adapted learners. *ArXiv.* 1896

Huajian Xin, Daya Guo, Zhihong Shao, Zhizhou Ren, 1897

Qihao Zhu, Bo Liu (Benjamin Liu), Chong Ruan, 1898

Wenda Li, and Xiaodan Liang. 2024a. **Deepseek-** 1899

prover: Advancing theorem proving in llms through 1900

large-scale synthetic data. *ArXiv*, abs/2405.14333. 1901

Huajian Xin, Z. Z. Ren, Junxiao Song, Zhihong 1902

Shao, Wanxia Zhao, Haocheng Wang, Bo Liu, Liyue 1903

Zhang, Xuan Lu, Qiushi Du, Wenjun Gao, Qihao 1904

Zhu, Dejian Yang, Zhibin Gou, Z. F. Wu, Fuli Luo, 1905

and Chong Ruan. 2024b. **Deepseek-prover-v1.5:** 1906

Harnessing proof assistant feedback for reinforce- 1907

ment learning and monte-carlo tree search. *arXiv.* 1908

Guangzhi Xiong, Eric Xie, Amir Hassan Shariatmadari, 1909

Sikun Guo, Stefan Bekiranov, and Aidong Zhang. 1910

2024a. Improving scientific hypothesis generation 1911

with knowledge grounded large language models. 1912

arXiv preprint arXiv:2411.02382. 1913

Siheng Xiong, Delin Chen, Qingyang Wu, Longxuan 1914

Yu, Qingzhen Liu, Dawei Li, Zhikai Chen, Xiaoze 1915

Liu, and Liangming Pan. 2024b. Improving causal 1916

reasoning in large language models: A survey. *arXiv* 1917

preprint arXiv:2410.16676. 1918

- 1919 Tianlong Xu, Yi-Fan Zhang, Zhendong Chu, Shen
1920 Wang, and Qingsong Wen. 2024a. Ai-driven virtual
1921 teacher for enhanced educational efficiency:
1922 Leveraging large pretrain models for autonomous
1923 error analysis and correction. *arXiv preprint*
1924 *arXiv:2409.09403*.
- 1925 Wanghan Xu, Yuhao Zhou, Yifan Zhou, Qinglong
1926 Cao, Shuo Li, Jia Bu, Bo Liu, Yixin Chen,
1927 Xuming He, Xiangyu Zhao, and 1 others. 2025.
1928 Probing scientific general intelligence of llms
1929 with scientist-aligned workflows. *arXiv preprint*
1930 *arXiv:2512.16969*.
- 1931 Yifan Xu, Xiao Liu, Xinghan Liu, Zhenyu Hou,
1932 Yueyan Li, Xiaohan Zhang, Zihan Wang, Aohan
1933 Zeng, Zhengxiao Du, Wenyi Zhao, Jie Tang, and
1934 Yuxiao Dong. 2024b. **Chatglm-math: Improving
1935 math problem-solving in large language models with
1936 a self-critique pipeline**. *ArXiv*, abs/2404.02893.
- 1937 Ziwei Xu, Sanjay Jain, and Mohan Kankanhalli.
1938 2024c. Hallucination is inevitable: An innate lim-
1939 itation of large language models. *arXiv preprint*
1940 *arXiv:2401.11817*.
- 1941 Yibo Yan and Joey Lee. 2024. Georeasoner: Reason-
1942 ing on geospatially grounded context for natural lan-
1943 guage understanding. In *Proceedings of the 33rd
1944 ACM International Conference on Information and
1945 Knowledge Management*, pages 4163–4167.
- 1946 Yibo Yan, Jiamin Su, Jianxiang He, Fangteng Fu,
1947 Xu Zheng, Yuanhuiyi Lyu, Kun Wang, Shen Wang,
1948 Qingsong Wen, and Xuming Hu. 2024a. A survey
1949 of mathematical reasoning in the era of multimodal
1950 large language model: Benchmark, method & chal-
1951 lenges. *arXiv preprint arXiv:2412.11936*.
- 1952 Yibo Yan, Shen Wang, Jiahao Huo, Hang Li, Boyan
1953 Li, Jiamin Su, Xiong Gao, Yi-Fan Zhang, Tianlong
1954 Xu, Zhendong Chu, and 1 others. 2024b. Errorradar:
1955 Benchmarking complex mathematical reasoning of
1956 multimodal large language models via error detec-
1957 tion. *arXiv preprint arXiv:2410.04509*.
- 1958 Yibo Yan, Shen Wang, Jiahao Huo, Philip S Yu, Xum-
1959 ing Hu, and Qingsong Wen. 2025. Mathagent:
1960 Leveraging a mixture-of-math-agent framework for
1961 real-world multimodal mathematical error detection.
1962 *arXiv preprint arXiv:2503.18132*.
- 1963 Yibo Yan, Haomin Wen, Siru Zhong, Wei Chen,
1964 Haodong Chen, Qingsong Wen, Roger Zimmer-
1965 mann, and Yuxuan Liang. 2024c. Urbanclip: Learn-
1966 ing text-enhanced urban region profiling with con-
1967 trastive language-image pretraining from the web.
1968 In *Proceedings of the ACM on Web Conference 2024*,
1969 pages 4006–4017.
- 1970 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng,
1971 Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan
1972 Li, Dayiheng Liu, Fei Huang, and 1 others.
1973 2024a. Qwen2 technical report. *arXiv preprint*
1974 *arXiv:2407.10671*.
- An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao,
Bowen Yu, Chengpeng Li, Dayiheng Liu, Jian-
hong Tu, Jingren Zhou, Junyang Lin, Keming Lu,
Mingfeng Xue, Runji Lin, Tianyu Liu, Xingzhang
Ren, and Zhenru Zhang. 2024b. **Qwen2.5-math
technical report: Toward mathematical expert model
via self-improvement**. *ArXiv*, abs/2409.12122.
- Liu Yang, Haihua Yang, Wenjun Cheng, Lei Lin,
Chenxia Li, Yifu Chen, Lunan Liu, Jianfei Pan, Tian-
wen Wei, Biye Li, Liang Zhao, Lijie Wang, Bo Zhu,
Guoliang Li, Xuejie Wu, Xilin Luo, and Rui Hu.
2023a. **Skymath: Technical report**. *arXiv preprint*
arXiv: 2310.16713.
- Yifei Yang, Runhan Shi, Z. Li, Shu Jiang, Bao-Liang
Lu, Yang Yang, and Hai Zhao. 2024c. **Batgpt-chem:
A foundation large model for retrosynthesis predic-
tion**. *ArXiv*, abs/2408.10285.
- Zhen Yang, Jinhao Chen, Zhengxiao Du, Wenmeng Yu,
Weihan Wang, Wenyi Hong, Zhihuan Jiang, Bin Xu,
Yuxiao Dong, and Jie Tang. 2024d. **Mathglm-vision:
Solving mathematical problems with multi-modal
large language model**. *ArXiv*, abs/2409.13729.
- Zhen Yang, Ming Ding, Qingsong Lv, Zhihuan Jiang,
Zehai He, Yuyi Guo, Jinfeng Bai, and Jie Tang.
2023b. Gpt can solve mathematical problems with-
out a calculator. *arXiv preprint arXiv:2309.03241*.
- Zonglin Yang, Wanhao Liu, Ben Gao, Yujie Liu,
Wei Li, Tong Xie, Lidong Bing, Wanli Ouyang,
Erik Cambria, and Dongzhan Zhou. 2025. Moose-
chem2: Exploring llm limits in fine-grained sci-
entific hypothesis discovery via hierarchical search.
arXiv preprint arXiv:2505.19209.
- Zonglin Yang, Wanhao Liu, Ben Gao, Tong Xie,
Yuqiang Li, Wanli Ouyang, Soujanya Poria, Erik
Cambria, and Dongzhan Zhou. 2024e. **Moose-
chem: Large language models for rediscovering
unseen chemistry scientific hypotheses**. *ArXiv*,
abs/2410.07076.
- Huanjin Yao, Jiaying Huang, Wenhao Wu, Jingyi
Zhang, Yibo Wang, Shunyu Liu, Yingjie Wang,
Yuxin Song, Haocheng Feng, Li Shen, and 1 others.
2024. Mulberry: Empowering mllm with o1-like
reasoning and reflection via collective monte carlo
tree search. *arXiv preprint arXiv:2412.18319*.
- Xinwu Ye, Chengfan Li, Siming Chen, Wei Wei, and
Robert Tang. 2025. **MMSciBench: Benchmarking
language models on Chinese multimodal scientific
problems**. In *Findings of the Association for Com-
putational Linguistics: ACL 2025*, pages 14621–
14663, Vienna, Austria. Association for Computa-
tional Linguistics.
- Huaiyuan Ying, Shuo Zhang, Linyang Li, Zhejian
Zhou, Yunfan Shao, Zhaoye Fei, Yichuan Ma, Ji-
awei Hong, Kuikun Liu, Ziyi Wang, Yudong Wang,
Zijian Wu, Shuaibin Li, Fengzhe Zhou, Hongwei
Liu, Songyang Zhang, Wenwei Zhang, Hang Yan,
Xipeng Qiu, and 3 others. 2024. **Internlm-math:**

2032	Open math large language models toward verifiable reasoning. <i>ArXiv</i> , abs/2402.06332.	2089
2033		2090
2034	Long Long Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T. Kwok, Zheng Li, Adrian Weller, and Weiyang Liu. 2023. Meta-math: Bootstrap your own mathematical questions for large language models . <i>ArXiv</i> , abs/2309.12284.	2091
2035		2092
2036		2093
2037		2094
2038		2095
2039	Sungduk Yu, Walter Hannah, Liran Peng, Jerry Lin, Mohamed Aziz Bhouri, Ritwik Gupta, Björn Lütjens, Justus C Will, Gunnar Behrens, Julius Busecke, and 1 others. 2024. Climsim: A large multi-scale dataset for hybrid physics-ml climate emulation. <i>Advances in Neural Information Processing Systems</i> , 36.	2096
2040		2097
2041		2098
2042		2099
2043		2100
2044		2101
2045		2102
2046	Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, and 1 others. 2024a. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 9556–9567.	2103
2047		2104
2048		2105
2049		2106
2050		2107
2051		2108
2052		2109
2053		2110
2054	Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wenhao Huang, Huan Sun, Yu Su, and Wenhua Chen. 2023. Mammoth: Building math generalist models through hybrid instruction tuning . <i>ArXiv</i> , abs/2309.05653.	2111
2055		2112
2056		2113
2057		2114
2058		2115
2059	Xiang Yue, Tianyu Zheng, Yuansheng Ni, Yubo Wang, Kai Zhang, Shengbang Tong, Yuxuan Sun, Botao Yu, Ge Zhang, Huan Sun, and 1 others. 2024b. Mmmu-pro: A more robust multi-discipline multimodal understanding benchmark. <i>arXiv preprint arXiv:2409.02813</i> .	2116
2060		2117
2061		2118
2062		2119
2063		2120
2064		2121
2065	Xiang Yue, Tuney Zheng, Ge Zhang, and Wenhua Chen. 2024c. Mammoth2: Scaling instructions from the web . <i>ArXiv</i> , abs/2405.03548.	2122
2066		2123
2067		2124
2068	Mohd Zaki, NM Anoop Krishnan, and 1 others. 2024. Mascca: investigating materials science knowledge of large language models. <i>Digital Discovery</i> , 3(2):313–327.	2125
2069		2126
2070		2127
2071		2128
2072	Zhiyuan Zeng, Qinyuan Cheng, Zhangyue Yin, Yunhua Zhou, and Xipeng Qiu. 2025. Revisiting the test-time scaling of o1-like models: Do they truly possess test-time scaling capabilities? <i>arXiv preprint arXiv:2502.12215</i> .	2129
2073		2130
2074		2131
2075		2132
2076		2133
2077	Di Zhang, Wei Liu, Qian Tan, Jingdan Chen, Hang Yan, Yuliang Yan, Jiatong Li, Weiran Huang, Xiangyu Yue, Dongzhan Zhou, Shufei Zhang, Mao Su, Han sen Zhong, Yuqiang Li, and Wanli Ouyang. 2024a. Chemllm: A chemical large language model . <i>ArXiv</i> , abs/2402.06852.	2134
2078		2135
2079		2136
2080		2137
2081		2138
2082		2139
2083	Di Zhang, Jianbo Wu, Jingdi Lei, Tong Che, Jiatong Li, Tong Xie, Xiaoshui Huang, Shufei Zhang, Marco Pavone, Yuqiang Li, and 1 others. 2024b. Llamaberry: Pairwise optimization for o1-like olympiad-level mathematical reasoning. <i>arXiv preprint arXiv:2410.02884</i> .	2140
2084		2141
2085		2142
2086		2143
2087		2144
2088		2145
	Ge Zhang, Xinrun Du, Bei Chen, Yiming Liang, Tongxu Luo, Tianyu Zheng, Kang Zhu, Yuyang Cheng, Chunpu Xu, Shuyue Guo, and 1 others. 2024c. Cmmmu: A chinese massive multi-discipline multimodal understanding benchmark. <i>arXiv preprint arXiv:2401.11944</i> .	2089
		2090
		2091
		2092
		2093
		2094
	Qiang Zhang, Keyang Ding, Tianwen Lyv, Xinda Wang, Qingyu Yin, Yiwen Zhang, Jing Yu, Yuhao Wang, Xiaotong Li, Zhuoyi Xiang, and 1 others. 2024d. Scientific large language models: A survey on biological & chemical domains. <i>arXiv preprint arXiv:2401.14656</i> .	2095
		2096
		2097
		2098
		2099
		2100
	Qingyang Zhang, Yake Wei, Zongbo Han, Huazhu Fu, Xi Peng, Cheng Deng, Qinghua Hu, Cai Xu, Jie Wen, Di Hu, and 1 others. 2024e. Multimodal fusion on low-quality data: A comprehensive survey. <i>arXiv preprint arXiv:2404.18947</i> .	2101
		2102
		2103
		2104
		2105
	Renrui Zhang, Dongzhi Jiang, Yichi Zhang, Haokun Lin, Ziyu Guo, Pengshuo Qiu, Aojun Zhou, Pan Lu, Kai-Wei Chang, Yu Qiao, and 1 others. 2025. Mathverse: Does your multi-modal llm truly see the diagrams in visual math problems? In <i>European Conference on Computer Vision</i> , pages 169–186. Springer.	2106
		2107
		2108
		2109
		2110
		2111
		2112
	Wenxuan Zhang, Mahani Aljunied, Chang Gao, Yew Ken Chia, and Lidong Bing. 2023. M3exam: A multilingual, multimodal, multilevel benchmark for examining large language models. <i>Advances in Neural Information Processing Systems</i> , 36:5484–5505.	2113
		2114
		2115
		2116
		2117
	Yu Zhang, Xiusi Chen, Bowen Jin, Sheng Wang, Shuiwang Ji, Wei Wang, and Jiawei Han. 2024f. A comprehensive survey of scientific large language models and their applications in scientific discovery. <i>arXiv preprint arXiv:2406.10833</i> .	2118
		2119
		2120
		2121
		2122
	Yuchen Zhang, Ratish Kumar Chandrakant Jha, Soumya Bharadwaj, Vatsal Sanjaykumar Thakkar, Adrienne Hoarfrost, and Jin Sun. 2024g. A benchmark dataset for multimodal prediction of enzymatic function coupling dna sequences and natural language. <i>arXiv preprint arXiv:2407.15888</i> .	2123
		2124
		2125
		2126
		2127
		2128
	Andrew Zhao, Daniel Huang, Quentin Xu, Matthieu Lin, Yong-Jin Liu, and Gao Huang. 2024a. Expel: Llm agents are experiential learners. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pages 19632–19642.	2129
		2130
		2131
		2132
		2133
	Wayne Xin Zhao, Kun Zhou, Zheng Gong, Beichen Zhang, Yuanhang Zhou, Jing Sha, Zhigang Chen, Shijin Wang, Cong Liu, and Ji rong Wen. 2022. Jiuzhang: A chinese pre-trained language model for mathematical problem understanding . <i>Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining</i> .	2134
		2135
		2136
		2137
		2138
		2139
		2140
	Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, and 1 others. 2023a. A survey of large language models. <i>arXiv preprint arXiv:2303.18223</i> .	2141
		2142
		2143
		2144
		2145

- Wayne Xin Zhao, Kun Zhou, Beichen Zhang, Zheng Gong, Zhipeng Chen, Yuanhang Zhou, Ji rong Wen, Jing Sha, Shijin Wang, Cong Liu, and Guoping Hu. 2023b. **Jiuzhang 2.0: A unified chinese pre-trained language model for multi-task mathematical problem solving**. *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*.
- Yilun Zhao, Lujing Xie, Haowei Zhang, Guo Gan, Yitao Long, Zhiyuan Hu, Tongyan Hu, Weiyuan Chen, Chuhan Li, Junyang Song, Zhijian Xu, Chengye Wang, Weifeng Pan, Ziyao Shanguan, Xiangru Tang, Zhenwen Liang, Yixin Liu, Chen Zhao, and Arman Cohan. 2025. **Mmvu: Measuring expert-level multi-discipline video understanding**. *Preprint*, arXiv:2501.12380.
- Zihan Zhao, Bo Chen, Jingpiao Li, Lu Chen, Liyang Wen, Pengyu Wang, Zichen Zhu, Danyang Zhang, Ziping Wan, Yansi Li, Zhongyang Dai, Xin Chen, and Kai Yu. 2024b. **Chemdfm-x: Towards large multimodal model for chemistry**. *ArXiv*, abs/2409.13194.
- Zihan Zhao, Da Ma, Lu Chen, Liangtai Sun, Zihao Li, Hongshen Xu, Zichen Zhu, Su Zhu, Shuai Fan, Guodong Shen, Xin Chen, and Kai Yu. 2024c. **Chemdfm: A large language foundation model for chemistry**. *ArXiv*.
- Kening Zheng, Junkai Chen, Yibo Yan, Xin Zou, and Xuming Hu. 2024. Reefknot: A comprehensive benchmark for relation hallucination evaluation, analysis and mitigation in multimodal large language models. *arXiv preprint arXiv:2408.09429*.
- Siru Zhong, Xixuan Hao, Yibo Yan, Ying Zhang, Yangqiu Song, and Yuxuan Liang. 2024a. Urban-cross: Enhancing satellite image-text retrieval with cross-domain adaptation. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 6307–6315.
- Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied, Weizhu Chen, and Nan Duan. 2024b. **AGIEval: A human-centric benchmark for evaluating foundation models**. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 2299–2314, Mexico City, Mexico. Association for Computational Linguistics.
- Guanyu Zhou, Yibo Yan, Xin Zou, Kun Wang, Aiwei Liu, and Xuming Hu. 2024a. Mitigating modality prior-induced hallucinations in multimodal large language models via deciphering attention causality. *arXiv preprint arXiv:2410.04780*.
- Kun Zhou, Beichen Zhang, Jiapeng Wang, Zhipeng Chen, Wayne Xin Zhao, Jing Sha, Zhichao Sheng, Shijin Wang, and Ji-Rong Wen. 2024b. **Jiuzhang3.0: Efficiently improving mathematical reasoning by training small data synthesis models**. *ArXiv*, abs/2405.14365.
- Xibin Zhou, Chenchen Han, Yingqi Zhang, Jin Su, Kai Zhuang, Shiyu Jiang, Zichen Yuan, Wei Zheng, Fengyuan Dai, Yuyang Zhou, Yuyang Tao, Dan Wu, and Fajie Yuan. 2025. **Decoding the molecular language of proteins with evolla**. *bioRxiv*.
- Wenwen Zhuang, Xin Huang, Xiantao Zhang, and Jin Zeng. 2024. Math-puma: Progressive upward multimodal alignment to enhance mathematical reasoning. *arXiv preprint arXiv:2408.08640*.
- Xin Zou, Yizhou Wang, Yibo Yan, Sirui Huang, Kening Zheng, Junkai Chen, Chang Tang, and Xuming Hu. 2024. Look twice before you answer: Memory-space visual retracing for hallucination mitigation in multimodal large language models. *arXiv preprint arXiv:2410.03577*.

2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229

2230
2231

2232
2233
2234
2235
2236
2237
2238
2239
2240

2241
2242
2243
2244
2245
2246
2247

2248
2249

2250
2251
2252
2253
2254
2255
2256

2257

2258
2259
2260
2261

Contents of Technical Appendices

A Clarification of the Position Scope	24
A.1 Novelty of Our Position	24
A.2 Unique Challenges for MLLMs in Scientific Reasoning vs. General Challenges	25
A.3 Significance and Examples of Reasoning across Scientific Domains .	26
A.4 Integration with Experimentation with Our Roadmap	26
A.5 Quantitative Analysis of Multimodal Components	26
B Formulation of Scientific Reasoning Task	27
C Four Phases of Research Roadmap	28
C.1 Phase 1: Broad Knowledge and Recognition	28
C.2 Phase 2: Analogical Reasoning and Generalization	29
C.3 Phase 3: Insightful Inference . . .	29
C.4 Phase 4: Creative Hypothesis Generation	29
C.5 Summary	29
D Data Differences among Four Domains	29
D.1 Mathematics: Structured Abstraction	29
D.2 Physics: Real-world Dynamics . .	29
D.3 Chemistry: Molecular and Symbolic	29
D.4 Biology: Conceptual Complexity .	30
D.5 Summary and Future Directions .	30
E Multi-Domain Scientific Reasoning Benchmark	30
F Details of Hallucinations in Scientific Reasoning	31
F.1 Different Types of Hallucinations in Scientific Reasoning	31
F.2 Existing Mitigation Strategies . .	33
F.3 Why Current Mitigation Methods Fall Short in Scientific Reasoning	33
G Multimodal Scientific (M)LLMs	34
H Discussion of More Modalities in Scientific Reasoning	34
H.1 Current Limitations and the Imperative for Richer Modalities	34

H.2 Technical Approaches for Broader Modality Integration	35	2262
H.3 Advanced Cross-Modal Interaction and Scientific Reasoning . . .	35	2264
H.4 Task-Driven Training Strategies and Future Directions	36	2266
I Discussion of Alternative Views	36	2268
I.1 Domain-Specific Models as a Superior Alternative	36	2269
I.2 Risks of Over-reliance on MLLMs	38	2270
J Clarification of LLM Usage	39	2271

Technical Appendices and Supplements

In this appendix, we first provide a clarification of our position’s scope (Appendix A), detailing its novelty, the unique challenges MLLMs face in scientific reasoning compared to general tasks, the significance of cross-domain reasoning, and the integration of experimentation with our proposed roadmap. Then, in Appendix B, we offer a formal formulation of the scientific reasoning task, while Appendix C elaborates on our four-stage research roadmap. Appendix D discusses data heterogeneity across key scientific domains, and Appendix E reviews relevant multi-domain scientific reasoning benchmarks, with Appendix G providing an overview of current scientific (M)LLMs. In Appendix F, we present an in-depth analysis of hallucinations in scientific reasoning, including their types, existing mitigation strategies, and why they often fall short in scientific contexts. Appendix H further explores the necessity and technical approaches for integrating richer modalities such as audio, video, and 3D data. Finally, Appendix I addresses alternative views and counterarguments, such as the preference for hyper-specialized models and concerns regarding MLLM reliability and explainability.

A Clarification of the Position Scope

A.1 Novelty of Our Position

This position paper presents a novel and compelling argument for the transformative potential of MLLMs in advancing the full spectrum of scientific reasoning, coupled with a unique, MLLM-centric four-stage roadmap towards this goal. While individual concepts such as applying AI to science or outlining stages of AI development have been discussed elsewhere, our work distinguishes itself through a specific synthesis, focus, and scope that addresses critical gaps in the current discourse:

1. Primacy of Multimodality for Comprehensive Scientific Reasoning: Our central thesis is that the inherent ability of MLLMs to integrate and reason over diverse data modalities is paramount for a holistic advancement in scientific reasoning. This contrasts sharply with existing literature that often focuses on unimodal approaches or narrower applications.

- For instance, comprehensive surveys on “Scientific LLMs” (Zhang et al., 2024f)

predominantly analyze text-based models and their applications in scientific discovery, largely overlooking the rich multimodal nature of scientific data and reasoning. Our position explicitly champions the multimodal perspective as indispensable.

- Similarly, while works like (Reddy and Shojaee, 2024) discuss generative AI for “scientific discovery,” their scope is primarily confined to this (albeit important) sub-task. Our position encompasses a broader range of scientific reasoning capabilities crucial for overall scientific advancement, including foundational knowledge acquisition, commonsense scientific question answering, analogical reasoning across disciplines, and insightful inference from complex data all areas where MLLMs offer unique advantages.

2. A Dedicated Four-Stage Roadmap for MLLM-Driven Scientific Reasoning Towards AGI: We propose a structured four-stage research roadmap (Broad Knowledge and Recognition, Analogical Reasoning and Generalization, Insightful Inference, and Creative Hypothesis Generation) that is specifically tailored to the evolving capabilities of MLLMs within *all scientific reasoning scenarios*, ultimately aiming towards Artificial General Intelligence (AGI). This roadmap is distinct from others in its MLLM-centricity and breadth:

- It moves beyond general discussions on generative AI’s potential in science (Morris, 2023), which, while insightful, did not offer a structured, MLLM-focused developmental pathway.
- It differs significantly from task-specific pipelines. For example, the highly valuable “AI Scientist” framework proposed by (Lu et al., 2024a) outlines a three-stage pipeline (idea generation → experiment iteration → paper write-up) primarily designed for *automated scientific discovery and specific generative outputs* like code and papers. Our roadmap is more encompassing: it situates such advanced discovery and generation capabilities (which we discuss in Section 5, including applications like digital teaching assistants)

2372 within its sophisticated later stages (partic- 2421
2373 ularly Stage 4). Crucially, our roadmap 2422
2374 also articulates the foundational and in- 2423
2375 termediate reasoning abilities (Stages 1- 2424
2376 3) that are prerequisite for achieving such 2425
2377 creative and autonomous scientific ex- 2426
2378 ploration with MLLMs. These earlier 2427
2379 stages, focusing on multimodal data in- 2428
2380 tegration, knowledge retrieval, contextual 2429
2381 understanding, and analogical generaliza- 2430
2382 tion, are essential for building robust and 2431
2383 reliable scientific reasoning systems. 2432

2384 In essence, the novelty of our position lies not 2433
2385 in addressing these themes in isolation, but in the 2434
2386 **holistic and integrated argument** that MLLMs, 2435
2387 by virtue of their multimodal capabilities, are 2436
2388 uniquely positioned to revolutionize scientific rea- 2437
2389 soning across its entire breadth. Furthermore, we 2438
2390 provide a **dedicated, MLLM-specific, and pro-** 2439
2391 **gressive developmental trajectory** towards real- 2440
2392 izing this vision. This distinct stance, we believe, 2441
2393 fills a crucial gap and warrants greater exposure 2442
2394 and discussion within the machine learning com- 2443
2395 munity, aligning perfectly with the objectives of a 2444
2396 Position Paper to highlight compelling arguments 2445
2397 that can shape future research directions. 2446

2398 **A.2 Unique Challenges for MLLMs in** 2447 2399 **Scientific Reasoning vs. General** 2448 2400 **Challenges** 2449

2401 While MLLMs face a spectrum of general chal- 2450
2402 lenges, their application to the domain of scien- 2451
2403 tific reasoning introduces distinct and often ampli- 2452
2404 fied difficulties. We identify four key areas where 2453
2405 the challenges for MLLMs in scientific reason- 2454
2406 ing diverge significantly from those encountered 2455
2407 in more general applications: 2456

2408 **1. Demand for Rigor, Precision, and Verifiabil-** 2457
2409 **ity:** In general tasks such as summarization or 2458
2410 creative writing, MLLMs may employ approxi- 2459
2411 mate reasoning and provide plausible-sounding 2460
2412 answers where occasional logical flaws or im- 2461
2413 precision might be acceptable. However, scien- 2462
2414 tific reasoning mandates an exceptionally 2463
2415 high standard. It requires strict logical consis- 2464
2416 tency, mathematical precision, and the ability 2465
2417 to produce verifiable step-by-step derivations. 2466
2418 The outputs must not only be correct but also 2467
2419 demonstrably so, adhering to the rigorous vali- 2468
2420 dation principles of the scientific method. 2469
2470

2. **Highly Specialized and Structured Multi-** 2421
2422 **modal Data:** General-purpose MLLMs are 2423
2424 typically designed to handle diverse, often un- 2425
2426 structured or semi-structured, multimodal data, 2427
2428 with alignment efforts often focusing on broad 2429
2430 semantic correspondence. In contrast, scien- 2431
2432 tific reasoning necessitates the processing and 2433
2434 deep understanding of highly structured, sym- 2435
2436 bolic, and specialized data modalities. This in- 2437
2438 cludes, but is not limited to, complex chemi- 2439
2440 cal formulas, genetic sequences, intricate dia- 2441
2442 grams with specific notational conventions, and 2443
2444 precise experimental data. Effective integra- 2444
2445 tion and reasoning over such specialized data 2445
2446 require more than general semantic understand- 2446
2447 ing; they demand the ability to parse and inter- 2447
2448 pret domain-specific syntax and semantics. 2448
2449

3. **The Critical Role of Causality and Mecha-** 2438
2439 **nistic Understanding:** While MLLMs often 2439
2440 excel at pattern recognition and identifying cor- 2440
2441 relations within large datasets—valuable capa- 2441
2442 bilities in many contexts—scientific reasoning 2442
2443 demands a more profound level of understand- 2443
2444 ing. A core objective of science is to move 2444
2445 beyond mere correlation to infer causality and 2445
2446 elucidate the underlying mechanisms that gov- 2446
2447 ern phenomena. MLLMs applied to science 2447
2448 must therefore develop capabilities to not just 2448
2449 describe *what* happens, but to reason about *why* 2449
2450 and *how* it happens, forming a crucial distinc- 2450
2451 tion from general pattern-matching tasks. 2451

4. **Nuanced Role and High Stakes of Halluci-** 2452
2453 **nation:** Hallucinations, or the generation of 2453
2454 factually incorrect or nonsensical information, 2454
2455 are generally undesirable in any MLLM appli- 2455
2456 cation as they lead to misinformation. How- 2456
2457 ever, in the context of scientific reasoning, fac- 2457
2458 tual hallucinations are particularly detrimental 2458
2459 due to the high stakes involved; inaccuracies 2459
2460 can misdirect research, lead to flawed 2460
2461 conclusions, and erode trust. Paradoxically, 2461
2462 there is also a unique potential role for *con-* 2462
2463 *trolled* generative capabilities—which might 2463
2464 share some characteristics with hallucination if 2464
2465 unconstrained—in the creative hypothesis gen- 2465
2466 eration phase. The challenge lies in fostering 2466
2467 this creative, “out-of-the-box” thinking while 2467
2468 strictly preventing uncontrolled factual inaccur- 2468
2469 acies, demanding a sophisticated balance not 2469
2470 typically required in general applications. 2470

A.3 Significance and Examples of Reasoning across Scientific Domains

The capacity to synthesize knowledge and reason across disparate scientific domains is fundamental to addressing complex, multifaceted global challenges and accelerating the pace of discovery. Many scientific frontiers lie at the intersection of disciplines, requiring an integrated understanding of phenomena that transcend traditional boundaries. MLLMs are uniquely positioned to facilitate this interdisciplinary synergy by their inherent ability to process, correlate, and reason over diverse data types—textual, visual, numerical, and symbolic—from various fields. This capability can unlock novel insights and solutions that might remain obscured when viewed through a single disciplinary lens. Below, we illustrate this potential with representative examples:

1. Drug Discovery and Precision Medicine:

MLLMs can integrate visual molecular structures, textual biomedical literature on biological pathways, and numerical patient genomic data. This allows for accelerated identification of promising drug candidates and prediction of patient-specific responses, paving the way for personalized medicine.

2. Materials Science and Engineering:

By analyzing material microstructure images (visual), crystal structure information (symbolic/visual), and textual data from scientific literature on synthesis and properties, MLLMs can rapidly predict characteristics of novel materials. This capability can guide experimental design and accelerate the discovery of materials with desired functionalities.

3. Climate Modeling and Earth System Science:

MLLMs can fuse satellite imagery (visual), numerical weather pattern data, and textual reports on atmospheric composition. This synthesis enables the construction of more comprehensive climate models, uncovering complex correlations and improving predictions of climate change impacts and extreme weather events.

A.4 Integration with Experimentation with Our Roadmap

Experimentation is undeniably central to scientific advancement. While our four-stage roadmap

(Broad Knowledge and Recognition, Analogical Reasoning, Insightful Inference, Creative Hypothesis Generation) focuses on the evolution of MLLM *reasoning capabilities*, it inherently supports the experimental process, which we view as a sophisticated workflow leveraging these abilities.

For instance, designing novel experiments draws upon Creative Hypothesis Generation (Stage 4) and Insightful Inference (Stage 3), while interpreting complex experimental data relies heavily on Insightful Inference. Our framework explicitly addresses this in Section 3.3 (“Simulation & Hypothesis Testing”), where MLLMs contribute to generating testable hypotheses, designing *in silico* experiments, and predicting outcomes—tasks demanding advanced inference and creative generation (Stages 3 & 4). The progression through our stages equips MLLMs to assist in experimental design (Stages 1 & 2), facilitate simulations (Stages 3 & 4), enhance multimodal data analysis (Stage 3), and support iterative refinement of hypotheses and experiments (Stages 3 & 4). Furthermore, our Discussion (Section 5) on the future outlook of MLLMs (input, model, output, environment) naturally aligns with the components of scientific experimentation, where MLLMs process experimental inputs, perform reasoning, generate outputs like protocols or analyzed results, and operate within the scientific environment.

In essence, the capacity to support and augment scientific experimentation is an emergent outcome as MLLMs advance through the foundational reasoning capabilities outlined in our roadmap, making them integral to the experimental lifecycle.

A.5 Quantitative Analysis of Multimodal Components

To substantiate our position on the challenges and opportunities for MLLMs in scientific reasoning, this section provides a quantitative analysis of how MLLMs process and rely on different modalities, with a particular focus on the visual component. The analysis draws upon empirical experiment from two comprehensive benchmarks, **SciVerse** (Guo et al., 2025b) and **MathVerse** (Zhang et al., 2025), which systematically deconstruct problem formulations to isolate and evaluate the models’ visual reasoning capabilities.

The central hypothesis investigated is that while MLLMs are designed to be multimodal, their per-

2570 formance in complex scientific and mathematical
2571 domains is disproportionately reliant on textual in-
2572 formation, often failing to genuinely "see" or rea-
2573 son from diagrams. This creates a significant per-
2574 formance gap when crucial information is shifted
2575 from text to the visual modality, mirroring many
2576 real-world scenarios.

Analysis from SciVerse: General Scientific Reasoning. The SciVerse benchmark evaluates MLLMs on scientific problems (Physics, Chemistry, Biology) by creating versions with varying degrees of visual dependency. The Knowledge-free version presents a standard problem with textual descriptions and a diagram. The Vision-rich version moves some conditional information from the text into the diagram, and the Vision-only version embeds the entire question within the visual input. As shown in Table 1, the performance of even top-tier models like GPT-4o degrades as the reliance on visual perception increases.

2591 The data reveals a consistent trend: when
2592 MLLMs are forced to extract information from di-
2593 agrams rather than text, their accuracy declines.
2594 GPT-4o’s accuracy drops by 5.0 percentage points,
2595 indicating that even state-of-the-art models struggle
2596 to interpret and integrate visually presented sci-
2597 entific conditions. This suggests that the visual en-
2598 coder and cross-modal fusion mechanisms in cur-
2599 rent MLLMs are a significant bottleneck for gen-
2600 eral scientific reasoning.

Analysis from MathVerse: Mathematical Reasoning. The MathVerse benchmark provides an even more granular analysis for mathematical problems, a domain where diagrams are often structured and indispensable. It progressively shifts information from text to vision across six versions. For this analysis, we focus on the versions that best illustrate the visual dependency: Text-dominant (full text), Text-lite (redundant text removed), Vision-intensive (implicit properties in vision), Vision-dominant (key values in vision), and Vision-only (all in vision).

2613 The results, as detailed in Table 2 and Fig-
2614 ure 6, are striking. The steady and signifi-
2615 cant performance degradation as we move from
2616 Text-dominant to Vision-only provides com-
2617 pelling evidence that MLLMs heavily rely on tex-
2618 tual "shortcuts." GPT-4V’s accuracy plummets by
2619 over 23 percentage points, demonstrating a crit-
2620 ical failure to parse essential conditions, prop-
2621 erties, and questions from mathematical diagrams.

2622 MathVerse further reveals that some models, such
2623 as Qwen-VL-Max, astonishingly perform better
2624 without the diagram input at all (+5.1%), treating
2625 the visual modality as a source of distraction rather
2626 than information.

2627 The quantitative data from both SciVerse and
2628 MathVerse converges on a clear conclusion: a crit-
2629 ical "vision-language gap" exists in the scientific
2630 reasoning capabilities of current MLLMs. Their
2631 strength lies predominantly in language process-
2632 ing, while their ability to perform deep, structured
2633 reasoning based on visual inputs especially the
2634 information-rich diagrams common in science and
2635 math remains underdeveloped. This quantitative
2636 analysis underscores the challenges highlighted in
2637 our position paper and reinforces the need for fu-
2638 ture research to focus on enhancing genuine mul-
2639 timodal understanding, moving beyond mere text-
2640 based reasoning with visual adornments.

2641 B Formulation of Scientific Reasoning 2642 Task

2643 Mathematically, the scientific reasoning process
2644 can be modeled as an optimization task over a mul-
2645 timodal knowledge graph $G = (V, E)$, where V
2646 denotes nodes (concepts, entities, or data points)
2647 and E represents edges (relationships or interac-
2648 tions). Let X_m represent the modality-specific in-
2649 put data, including mandatory modalities such as
2650 textual descriptions X_t and visual representations
2651 X_v , as well as optional modalities like numerical
2652 data X_n or other specialized inputs X_o . The goal
2653 is to predict or infer target outputs Y based on a
2654 reasoning function f_θ , parameterized by θ , over G :

$$2655 Y = f_\theta(G, X_m) = f_\theta(V, E, X_t, X_v, X_n, X_o).$$

2656 **Mathematics:** Consider solving a geometry
2657 problem where X_t provides the problem descrip-
2658 tion, X_v includes the geometric diagram, and X_n
2659 optionally contains measurements or coordinates.
2660 The reasoning function f_θ integrates these inputs
2661 to infer the solution, such as identifying the area
2662 of a triangle.

2663 **Physics:** In a physics experiment, X_t describes
2664 the theoretical background, X_v presents the ex-
2665 perimental setup image, and X_n includes sen-
2666 sor data such as velocity or temperature measure-
2667 ments. The function f_θ predicts the outcome or
2668 validates a hypothesis.

2669 **Chemistry:** A multimodal analysis of a chem-
2670 ical reaction could include X_t for the reaction

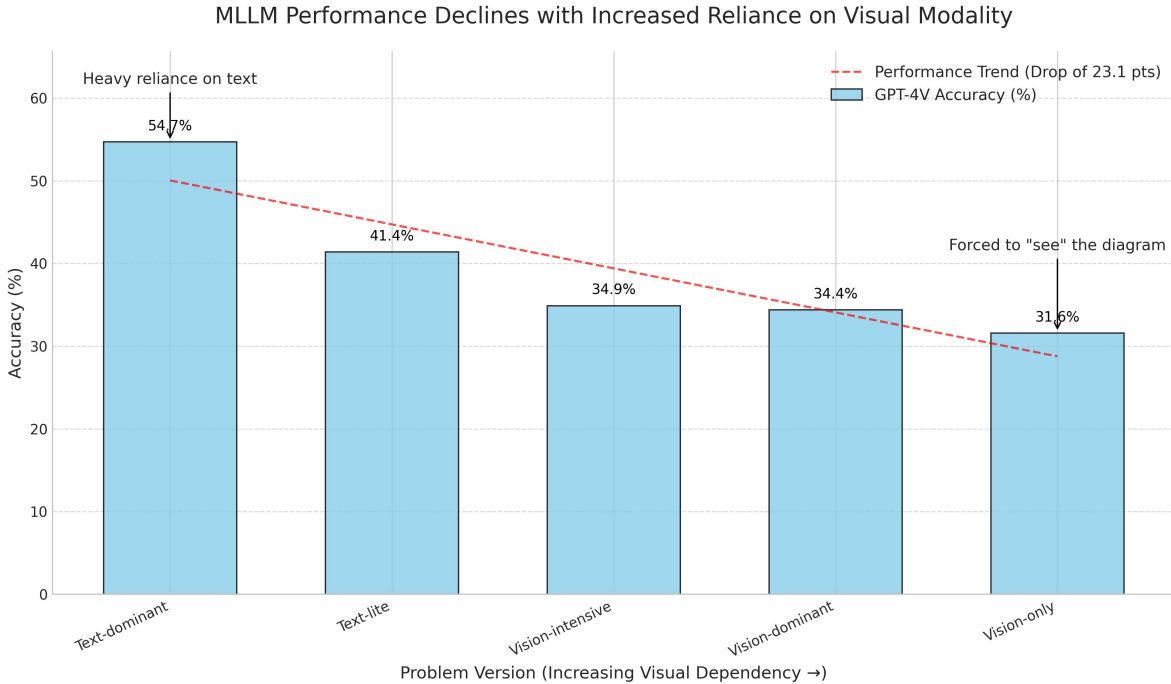


Figure 6: MLLM performance declines with increased reliance on visual modality from MathVerse benchmark.

Table 1: Accuracy of MLLMs on SciVerse versions with increasing visual dependency. The performance drop from Knowledge-free to Vision-only highlights the challenge of visual information processing and reasoning.

Model	Knowledge-free (Acc %)	Vision-rich (Acc %)	Vision-only (Acc %)	Perf. Drop (KF to VO)
GPT-4o	55.2	54.4	50.2	-5.0 pts
LLaVA-OneVision (7B)	47.2	46.5	41.9	-5.3 pts

Table 2: Accuracy of MLLMs across MathVerse versions, showing a progressive decline as problems become more vision-reliant.

Model	Text-dom (%)	Text-lite (%)	Vision-int (%)	Vision-dom (%)	Vision-only (%)	Overall Drop
GPT-4V	54.7	41.4	34.9	34.4	31.6	-23.1 pts
InternLM-XC2 (7B)	22.3	17.0	15.7	16.4	11.0	-11.3 pts

mechanism, X_v for molecular structure visualizations, and X_o for spectroscopic data. The model predicts reaction yield or product properties.

Biology: When studying gene expression, X_t describes the biological context, X_v contains microscope images, and X_n optionally includes numerical gene expression levels. The reasoning function predicts gene interactions or cellular behavior.

These examples illustrate how MLLMs can handle diverse inputs, integrating mandatory and optional modalities to perform complex scientific reasoning tasks.

C Four Phases of Research Roadmap

MLLMs have seen significant advancements in recent years, positioning their reasoning capabilities

as a pivotal element on the pathway to achieving Artificial General Intelligence (AGI) (Wang et al., 2024g; Yan et al., 2024a; Sun et al., 2023; Jin et al., 2024a; Yan et al., 2024c; Wei et al., 2024). However, realizing AGI requires navigating a structured roadmap characterized by progressively complex reasoning tasks. This section outlines the Four Phases of Research Roadmap, each delineated by unique data requirements, reasoning mechanisms, generalization abilities, and impact. Table 3 provides a comparative overview of these phases.

C.1 Phase 1: Broad Knowledge and Recognition

In this phase, MLLMs prioritize high diversity and low specificity in data and knowledge. Tasks involve integrating vast and diverse datasets, empha-

sizing retrieval and alignment mechanisms. For example, MLLMs excel in synthesizing encyclopedic knowledge, aligning visual inputs (*e.g.*, diagrams) with textual descriptions, and providing domain-specific insights. The generalization is limited to specific domains, resulting in low-impact applications like data retrieval and integration.

C.2 Phase 2: Analogical Reasoning and Generalization

This stage emphasizes medium diversity and specificity in data, leveraging contextual understanding to enable relational and analogical reasoning. For instance, MLLMs may draw analogies between chemical reaction pathways and electrical circuits, enhancing interdisciplinary insights. The models exhibit medium-level generalization across domains, allowing for moderate complexity tasks like explaining scientific phenomena using cross-domain analogies.

C.3 Phase 3: Insightful Inference

As reasoning tasks grow in complexity, MLLMs in this phase focus on low diversity but high specificity datasets. Predictive reasoning mechanisms are central, enabling context-specific inferences. For example, a model might predict the behavior of a physical system under certain constraints or optimize complex processes like material design. These capabilities lead to high-impact applications, including scientific optimization and inferential problem-solving.

C.4 Phase 4: Creative Hypothesis Generation

The final phase demands both high diversity and high creativity in data. Generative reasoning mechanisms empower MLLMs to propose innovative solutions and simulate hypotheses. For instance, models might design novel molecules for drug discovery or hypothesize ecological models for sustainable ecosystems. This phase represents very high-impact applications, bridging the gap between scientific discovery and innovation.

C.5 Summary

In summary, the Four Phases of Research Roadmap reflect the increasing complexity and potential of MLLMs in scientific reasoning. While current MLLMs have demonstrated impressive capabilities, they remain far from achieving AGI (Mumuni and Mumuni, 2025; Wang, 2024; Feng

et al., 2024; Fei et al., 2022). The community must continue to explore advanced reasoning abilities, fostering collaboration and innovation along this roadmap to address the challenges of AGI-driven scientific reasoning.

D Data Differences among Four Domains

The reasoning capabilities of MLLMs make them particularly well-suited for processing heterogeneous multimodal data (Pattnayak et al., 2024; Li et al., 2024e,i). By integrating and analyzing diverse data formats such as text, images, and structured information, MLLMs enable more comprehensive insights across various scientific disciplines. This section highlights the unique data characteristics of four domains: mathematics, physics, chemistry, and biology, as summarized in Table 4.

D.1 Mathematics: Structured Abstraction

Mathematical data is characterized by its symbolic equations, graphs, and geometric figures. These data types are highly structured and abstract, requiring precise interpretation and manipulation. Visual features such as coordinate axes and geometric shapes often complement formal textual elements like equations and proofs. The integration of these modalities enables MLLMs to solve complex mathematical problems and support theorem proving (Chern et al., 2023; Yan et al., 2024a; Azerbayev et al., 2023; Li et al., 2024i).

D.2 Physics: Real-world Dynamics

Physics data encompasses diagrams (*e.g.*, vector and circuit diagrams), graphs, and descriptions of real-world phenomena. Its data structure reflects system dynamics and real-world applications, combining descriptive text with visual representations like force vectors or particle motion. MLLMs leverage these multimodal inputs to model physical systems and predict outcomes under varying conditions (Jaiswal et al., 2024; Barman et al., 2025; Yu et al., 2024).

D.3 Chemistry: Molecular and Symbolic

Chemistry relies heavily on molecular structures, reaction pathways, and the periodic table. The data is both symbolic and structural, involving visual elements such as 3D molecular models and reaction schemes. Textual features include chemical equations and reaction mechanisms. MLLMs

Table 3: Comparison across four stages along key dimensions.

Dimension	Broad Knowledge and Recognition	Analogical Reasoning and Generalization	Insightful Inference	Creative Hypothesis Generation
Data and Knowledge Requirements	High diversity, low specificity	Medium diversity, medium specificity	Low diversity, high specificity	High diversity, high creativity
Reasoning Mechanisms	Low complexity (retrieval, alignment)	Medium complexity (relational, analogical)	High complexity (predictive)	Very high complexity (generative)
Model Generalization	Low (domain-specific)	Medium (cross-domain)	High (context-specific)	Very high (innovative)
Applications and Impact	Low impact (retrieval, integration)	Medium impact (interdisciplinary insights)	High impact (optimization, inference)	Very high impact (discovery, innovation)

facilitate understanding chemical interactions and even predicting new compounds or reaction outcomes (Alampara et al.; Miret and Krishnan, 2024; Mirza et al., 2024; Ramos et al., 2025; M. Bran et al., 2024).

D.4 Biology: Conceptual Complexity

Biological data is diverse, covering anatomical diagrams, cellular structures, and ecological models. Its structure is conceptual and often involves biological and anatomical representations. Visual inputs like cell diagrams and ecosystem models are paired with textual descriptions of processes and interactions. MLLMs support tasks such as identifying biological patterns and simulating ecological dynamics (Kraus et al., 2024; Luu and Buehler, 2023; Huang et al., 2024a; Wang et al., 2023a; Liu et al., 2023c; Tang et al., 2023).

D.5 Summary and Future Directions

In summary, MLLMs demonstrate strong reasoning capabilities across mathematics, physics, chemistry, and biology by integrating diverse multimodal data. However, their potential is not confined to these four disciplines. Future research should explore reasoning capabilities in broader domains such as geospatial analysis (Roberts et al., 2024; Mai et al., 2024; Hao et al., 2024; Zhong et al., 2024a) and coding (Li et al., 2024g; Guo et al., 2024a; Di et al., 2024; Hui et al., 2024), which demand advanced generalization abilities. Such endeavors are crucial for achieving the comprehensive reasoning capabilities envisioned in the AGI roadmap.

E Multi-Domain Scientific Reasoning Benchmark

The emergence of multi-domain scientific reasoning benchmarks has played a pivotal role in advancing AI models’ ability to reason across di-

verse scientific domains. These benchmarks vary in their focus on multimodal integration, educational levels, and the comprehensiveness of domain coverage, as summarized in Table 5.

Comprehensive Domain Coverage: Several benchmarks, such as **MMMU-Pro**, **CMMMU**, and **SciEval**, provide extensive domain coverage, including mathematics, physics, chemistry, and biology. These benchmarks target diverse educational contexts, ranging from primary education to PhD-level tasks, ensuring broad applicability. Notably, **MMMU-Pro** and **SciBench** have become instrumental for college-level evaluation, while **CMMMU** extends its scope to younger learners, addressing a critical gap in foundational scientific reasoning.

Multimodal Reasoning: The integration of multimodal capabilities has become a defining feature of contemporary benchmarks. Over 60% of the surveyed benchmarks, such as **EXAMS-V**, **SciBench**, and **ScienceQA**, incorporate multimodal tasks that involve textual, visual, and symbolic reasoning. These tasks reflect real-world problem-solving scenarios where multiple modalities interact, making them essential for evaluating the holistic reasoning capabilities of advanced AI models.

Specialized vs. General Benchmarks: While benchmarks like **OlympiadBench** and **OlympicArena** are tailored to high-stakes competitions, their narrow focus on specific tasks, such as Olympiad-level challenges, limits their generalizability. Conversely, resources like **AGIEval** and **SciEval** aim to provide a broader evaluation of scientific reasoning, covering multiple domains across various educational levels. However, their lack of multimodal integration highlights the need for more versatile benchmarks.

Future Directions: The future of multi-domain scientific reasoning benchmarks lies in developing

Table 4: Comparison of data features among four domains within the scope of this position paper.

Feature/Domain	Mathematics	Physics	Chemistry	Biology
Data Types	Symbolic equations, graphs, geometric figures	Diagrams (vector, circuit), graphs, real-world phenomena	Molecular structures, reaction pathways, periodic table	Anatomical diagrams, cellular structures, ecological models
Data Structure	Abstract, highly structured, formal notation	Real-world applications, system dynamics	Symbolic and molecular, structural	Conceptual, biological, and anatomical
Visual Features	Coordinate axes, geometric shapes	Diagrams e.g., force vectors, particle motion	3D molecular models, reaction schemes	Photos, cell diagrams, ecosystem models
Textural Features	Equations, proofs, formal definitions	Descriptive, experimental setups, theories	Chemical equations, reaction mechanisms	Biological processes, ecological interactions

unified frameworks that combine multimodal reasoning, domain comprehensiveness, and adaptability across educational levels. Incorporating elements like interactive feedback systems and collaborative reasoning tasks will bridge the gap between theoretical evaluation and practical applications, fostering the next generation of scientific reasoning models.

F Details of Hallucinations in Scientific Reasoning

Since a position paper must balance scope breadth and technical depth, we did not extensively elaborate on each specific challenge in the main text. However, we provide a more detailed discussion on hallucinations here, as they present a nuanced challenge particularly critical to the advancement of MLLMs in scientific reasoning. While often detrimental, understanding and managing hallucinations is key, especially considering the aspirational "Creative Hypothesis Generation" stage of our roadmap.

Recent studies on multimodal reasoning hallucination, such as MIRAGE (Dong et al., 2025), reveal complex and nuanced patterns that challenge one-size-fits-all mitigation strategies. As illustrated in Figure 7, the distribution of hallucination types is strongly correlated with the nature of the scientific task; for instance, logical hallucinations are pervasive across most question types, whereas factuality hallucinations are more prominent in knowledge-intensive domains like statistics, and spatial hallucinations are uniquely problematic in geometry and spatial tasks. Furthermore, Figure 8 demonstrates that while increasing model size can effectively reduce logical and fabrication hallucinations, it provides only marginal improvements for spatial hallucinations, indicating that this type of error stems from a deeper visual reasoning deficit that cannot be solved by simple

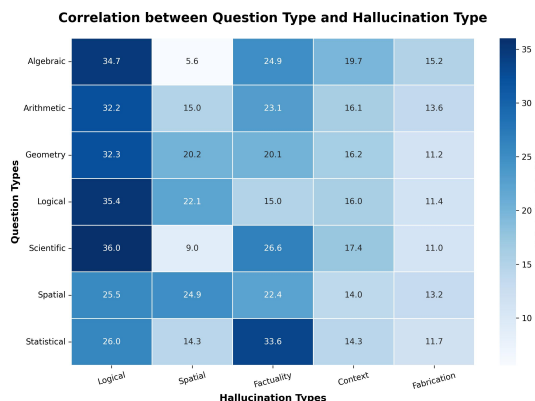


Figure 7: Eight prospects for the future of MLLMs in the field of multimodal scientific reasoning

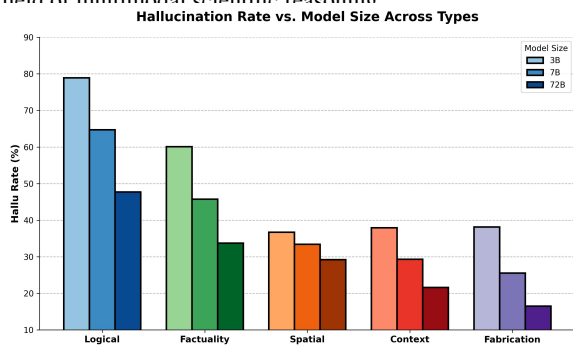


Figure 8: Eight prospects for the future of MLLMs in the field of multimodal scientific reasoning.

scaling. These findings underscore that the hallucination categories defined for multimodal scientific reasoning, which isolate reasoning failures from perceptual errors, are distinct and behave differently from general-purpose hallucination concepts. Their unique characteristics and resistance to conventional scaling solutions cannot be fully captured by broader definitions, thus necessitating the more targeted, fine-grained hallucination classification.

F.1 Different Types of Hallucinations in Scientific Reasoning

Hallucinations in MLLMs, broadly defined as the generation of information not grounded in the in-

Table 5: Comparison of multi-domain scientific reasoning benchmarks.

Paper	Organization	Venue	Multimodal	Education Level	Domain(s)			
					Math	Physics	Chemistry	Biology
MMMU-Pro (Yue et al., 2024b)	CMU	ACL'25	✓	College	✓	✓	✓	✓
SciVerse (Guo et al., 2025b)	CUHK	ACL Findings'25	✓	General	✓	✓	✓	✓
MMSciBench (Ye et al., 2025)	Yale University	ACL Findings'25	✓	High	✓	✓	✓	✓
EMMA (Hao et al., 2025)	UESTC	ICML'25	✓	General	✓	✓	✓	✓
LLM-SRBench (Shojaee et al., 2025)	Virginia Tech	ICML'25	✓	General	✓	✓	✓	✓
CURIE (Cui et al., 2025)	Google	ICLR'25	✓	General	✓	✓	✓	✓
ScienceAgentBench (Chen et al., 2024c)	OSU	ICLR'25	✓	General	✓	✓	✓	✓
Scientist-Bench (Tang et al., 2025)	HKU	NeurIPS'25	✓	General	✓	✓	✓	✓
MMVU (Zhao et al., 2025)	Yale University	CVPR'25	✓	College	✓	✓	✓	✓
SGL-Bench (Xu et al., 2025)	Shanghai AI Lab	arXiv'25	✓	General	✓	✓	✓	✓
SDE (Song et al., 2025)	Deep Principle	arXiv'25	✓	General	✓	✓	✓	✓
SciBench (Wang et al., 2024f)	UCLA	ICML'24	✓	College	✓	✓	✓	✓
MMMU (Yue et al., 2024a)	University of Waterloo & Ohio State University	CVPR'24	✓	College	✓	✓	✓	✓
EXAMS-V (Das et al., 2024)	MBZUAI	ACL'24	✓	General	✓	✓	✓	✓
ArxivQA/Cap (Li et al., 2024h)	HKU	ACL'24	✓	General	✓	✓	✓	✓
OlympiadBench (He et al., 2024a)	Tsinghua University	ACL'24	✓	Competition	✓	✓	✓	✓
ScemQA (Liang et al., 2024a)	University of Notre Dame	ACL Short'24	✓	College	✓	✓	✓	✓
OlympicArena (Huang et al., 2024b)	Shanghai Jiaotong University	NeurIPS'24	✓	Competition	✓	✓	✓	✓
SciEval (Sun et al., 2024b)	Shanghai Jiaotong University	AAAI'24	✓	General	✓	✓	✓	✓
CMMU (He et al., 2024c)	Beijing Academy of Artificial Intelligence	IJCAI'24	✓	Primary/Middle/High	✓	✓	✓	✓
AGIEval (Zhong et al., 2024b)	Microsoft	NAACL'24	✓	High/College	✓	✓	✓	✓
CMMLU (Li et al., 2024d)	MBZUAI	ACL Findings'24	✓	High/College	✓	✓	✓	✓
MMWorld (He et al., 2024b)	UCSC	arXiv'24	✓	General	✓	✓	✓	✓
CMMMU (Zhang et al., 2024c)	HKUST	arXiv'24	✓	College	✓	✓	✓	✓
M4U (Wang et al., 2024a)	Chinese Academy of Sciences	arXiv'24	✓	College	✓	✓	✓	✓
MMSci (Li et al., 2024m)	UCSB	arXiv'24	✓	PhD	✓	✓	✓	✓
LitQA2 (Skarlinski et al., 2024)	FutureHouse Inc.	arXiv'24	✓	General	✓	✓	✓	✓
VisScience (Jiang et al., 2024b)	Tsinghua	arXiv'24	✓	Primary/Middle/High	✓	✓	✓	✓
JEEBench (Arora et al., 2023)	Microsoft & UC Berkeley	EMNLP'23	✓	General	✓	✓	✓	✓
M3Exam (Zhang et al., 2023)	Alibaba	NeurIPS'23	✓	Primary/Middle/High	✓	✓	✓	✓
LitQA (Lála et al., 2023)	FutureHouse Inc.	arXiv'23	✓	General	✓	✓	✓	✓
ScienceQA (Lu et al., 2022a)	UCLA	NeurIPS'22	✓	Primary/High	✓	✓	✓	✓
IconQA (Lu et al., 2021b)	UCLA	NeurIPS'21	✓	General	✓	✓	✓	✓
MMLU (Hendrycks et al., 2021)	UC Berkeley	ICLR'21	✓	General	✓	✓	✓	✓

put data or established knowledge, manifest in several ways pertinent to scientific reasoning:

- **Factual Hallucination:** This is the most straightforward type, where the MLLM generates incorrect or non-existent facts, such as an incorrect physical constant, a misremembered chemical property, a wrong date for a discovery, or a fabricated citation. In science, where precision is paramount, factual hallucinations can lead to fundamentally flawed conclusions.
- **Conceptual Hallucination:** This involves the misrepresentation or misapplication of scientific concepts, theories, or principles. An MLLM might, for example, confuse correlation with causation in a biological context, misapply Newton's laws in a scenario where relativistic effects are dominant, or incorrectly define a mathematical theorem's conditions.
- **Relational Hallucination:** This occurs when an MLLM fabricates or misrepresents relationships between entities, variables, or concepts. It might invent a non-existent interaction between

two proteins, incorrectly link a gene to a disease without evidence, or propose a flawed causal chain in a physical process.

- **Explanatory Hallucination:** Here, the MLLM generates explanations for phenomena that are either logically unsound, inconsistent with established theories, or based on fabricated intermediate steps, even if the final conclusion appears plausible. This is particularly dangerous as it can create a false sense of understanding. For instance, an MLLM might "derive" a correct mathematical formula through a series of incorrect algebraic manipulations.
- **Exploratory Hallucination:** This is a more nuanced category. It refers to the generation of novel, plausible-sounding hypotheses or ideas that are not yet validated or directly supported by existing data. While technically a "hallucination" (as it's not grounded in current certainties), this type can be beneficial in the "Creative Hypothesis Generation" stage if properly managed and clearly identified as speculative. The risk

lies in presenting such exploratory outputs as established facts.

F.2 Existing Mitigation Strategies

Several strategies have been developed to mitigate hallucinations in general MLLM applications, which can be adapted, with limitations, to scientific reasoning:

- **Knowledge Retrieval Augmentation (RAG):**

This involves grounding the MLLM's responses by retrieving relevant information from external, trusted knowledge bases (*e.g.*, scientific databases, curated literature). The model then uses this retrieved context to formulate its answer, reducing reliance on its parametric memory which can be a source of hallucinations.

- **Consistency Checks:**

- *Self-consistency*: Generating multiple reasoning paths or answers for the same query and selecting the most consistent or frequently occurring one.
- *Cross-modal consistency*: For MLLMs, ensuring that information derived from text aligns with information from images (*e.g.*, a diagram in a physics problem should match the textual description).
- *Logical consistency*: Evaluating the logical coherence of the generated reasoning steps.

- **Process-Based Supervision (and Process Reward Models - PRMs):**

Instead of only rewarding the correctness of the final answer, these methods involve supervising or rewarding the intermediate steps of the reasoning process. This encourages the model to follow valid inferential pathways, making it less likely to "jump" to conclusions through hallucinated steps.

- **Human Feedback (RLHF/RLAIF):**

Reinforcement Learning from Human Feedback (and its variants like RLAIF) involves training the MLLM based on human ratings of its outputs. Humans can penalize hallucinated content, guiding the model towards more factual and reliable responses.

F.3 Why Current Mitigation Methods Fall Short in Scientific Reasoning

Despite their utility, current hallucination mitigation strategies face significant hurdles when applied to the rigorous demands of scientific reasoning:

- **Complexity and Nuance of Scientific Knowledge:**

- Scientific knowledge is vast, deeply interconnected, and often highly nuanced or domain-specific. RAG systems might struggle to retrieve the *exact* precise piece of information needed or may misinterpret the context of retrieved documents. The sheer volume can also lead to retrieving conflicting information.
- Many scientific concepts are abstract and symbolic (*e.g.*, in quantum mechanics or advanced mathematics), making it difficult for models to truly "understand" and ground them, even with retrieved text.

- **The Need for Exploration and Controlled Creativity:**

- Aggressive hallucination mitigation, designed to enforce strict factuality, can inadvertently stifle the model's ability to engage in exploratory reasoning or generate novel hypotheses (critical for Stage 4). The challenge is to allow for "beneficial" exploratory generation while suppressing detrimental factual or conceptual hallucinations. A balance current methods struggle to achieve.

- **Lack of Granular Explainability and Verifiability:**

- Even if a mitigation technique reduces overt hallucinations, the internal reasoning process of MLLMs often remains a black box. For scientific applications, it's not enough for an answer to be correct; the reasoning must be transparent, verifiable, and align with established scientific methods. Current methods don't inherently guarantee this level of scrutiny.
- Error propagation in multi-step scientific reasoning is a major issue. A subtle, uncaught hallucination in an early step can invalidate the entire chain, and PRMs might not be sufficiently fine-grained to catch all such nuanced errors in complex scientific domains.

- **Data Limitations for Scientific Domains:**

- High-quality, large-scale, and *verified* multimodal scientific datasets specifically cu-

rated for training MLLMs to avoid domain-specific hallucinations are scarce.

- Training effective PRMs or RLAIIF systems for science requires extensive, expert-annotated data on correct reasoning processes, which is costly and time-consuming to produce across diverse scientific fields.
- Models trained on existing scientific literature might inherit biases or outdated information present in that corpus, leading to "hallucinations" relative to the current state of knowledge.

In conclusion, while existing mitigation techniques provide a starting point, the unique demands of scientific reasoning—its need for precision, depth, verifiability, and controlled creativity—necessitate the development of more sophisticated, domain-aware, and interpretable approaches to manage and understand hallucinations in MLLMs.

G Multimodal Scientific (M)LLMs

Scientific (M)LLMs represent an emerging class of models that integrate multiple modalities to tackle complex problems across various scientific disciplines, as summarized in Table 6. These models leverage massive datasets combining text, images, graphs, and other forms of scientific data to provide deeper insights and advanced reasoning capabilities. Unlike traditional models, which typically focus on either textual or visual information, scientific MLLMs are designed to harmonize and synthesize diverse sources of scientific knowledge, enabling enhanced performance in tasks such as mathematical problem solving, chemical reaction prediction, biological analysis, and physical simulations (Morris, 2023; Pei et al., 2024; Reddy and Shojaee, 2024; Zhang et al., 2024f).

H Discussion of More Modalities in Scientific Reasoning

While current MLLMs have demonstrated remarkable progress, their sensory input is predominantly confined to text and static images. This limitation curtails their ability to fully comprehend and reason about the multifaceted nature of scientific phenomena, which often unfold dynamically, possess intricate three-dimensional structures, or generate non-visual sensory data. To significantly advance scientific reasoning, MLLMs must evolve to integrate a richer spectrum of modalities, including

audio, video, 3D structural data, and other sensor-derived information. This expansion is not merely about processing more data types but about enabling a more holistic and nuanced understanding essential for complex scientific inquiry.

H.1 Current Limitations and the Imperative for Richer Modalities

The current text-image paradigm in MLLMs, while powerful, presents several inherent limitations when applied to the diverse needs of scientific reasoning:

- **Restricted Modality Spectrum:** Most MLLMs excel with visual and textual data but struggle to incorporate or meaningfully reason over audio (*e.g.*, sonification of data, acoustic signatures of experiments), video (*e.g.*, dynamic processes, experimental procedures), and 3D data (*e.g.*, molecular configurations, geological strata). This restricts their ability to capture the full context of many scientific observations.
- **Inadequate for Dynamic Phenomena:** Scientific processes are often dynamic. For instance, observing a chemical reaction, the growth of a biological culture, or the time-evolution of a physical system requires processing temporal information embedded in video or sequential sensor readings. Static images and text alone cannot capture these crucial dynamic aspects.
- **Insufficient for Spatial Complexity:** Many scientific domains, such as chemistry (molecular structures), biology (protein folding, anatomical systems), materials science (crystal lattices), and earth sciences (geophysical models), fundamentally rely on understanding complex 3D spatial relationships. Current MLLMs often lack the sophisticated 3D geometric reasoning capabilities required.
- **Limited Cross-Modal Association for Complex Reasoning:** Scientific reasoning frequently involves drawing inferences across multiple, diverse data streams. For example, correlating a change in a visual indicator (image/video) with a specific sound (audio) and a corresponding data log (text/tabular) to deduce a causal link in an experiment demands finer-grained and more diverse modality interactions than currently supported.

- 3165 • **Domain-Specific Knowledge Integration** 3212
3166 **Challenges:** Effectively interpreting special- 3213
3167 ized modalities (*e.g.*, spectrographic data, 3214
3168 sonified outputs) necessitates not just raw 3215
3169 data processing but also deep integration
3170 with domain-specific ontologies and knowl-
3171 edge bases to contextualize the sensory input
3172 accurately.

3173 H.2 Technical Approaches for Broader 3216 3174 Modality Integration

3175 Addressing these limitations requires developing 3217
3176 robust technical approaches for encoding and 3218
3177 integrating a wider array of modalities, tailored 3219
3178 to their unique characteristics and scientific rele- 3220
3179 vance: 3221

3180 • Audio Processing: 3222

- 3181 – *Feature Extraction:* Leveraging estab- 3223
3182 lished techniques like spectrograms or Mel- 3224
3183 Frequency Cepstral Coefficients (MFCCs) 3225
3184 to convert raw audio signals from exper- 3226
3185 iments (*e.g.*, equipment sounds, material 3227
3186 stress responses) or sonified data into fea- 3228
3187 ture representations suitable for neural net- 3229
3188 works. 3230
3189 – *Model Adaptation:* Developing specialized 3231
3190 audio encoders or adapting existing ones, 3232
3191 and aligning their latent spaces with those 3233
3192 of text and image encoders through cross- 3234
3193 modal attention mechanisms or joint em- 3235
3194 bedding strategies. This allows, for in- 3236
3195 stance, an MLLM to associate the sound 3237
3196 of a failing component with its visual de- 3238
3197 pication and textual description of failure 3239
3198 modes. 3240

3199 • Video Understanding: 3241

- 3200 – *Spatiotemporal Modeling:* Employing 3242
3201 3D Convolutional Neural Networks (3D 3243
3202 CNNs), Video Transformers, or other 3244
3203 video-language models to capture both spa- 3245
3204 tial features within frames and temporal dy- 3246
3205 namics across frames. This is critical for 3247
3206 analyzing experimental procedures, observ- 3248
3207 ing cellular motility, or tracking particle tra- 3249
3208 jectories. 3250
3209 – *Keyframe and Event Detection:* Implement- 3251
3210 ing mechanisms to automatically identify 3252
3211 and focus on critical events or keyframes 3253

3212 within a scientific video (*e.g.*, phase tran- 3213
3214 sition, initiation of a reaction), thereby re- 3215
3216 ducing computational load and highlight- 3217
3218 ing salient information for reasoning. 3219

3220 • 3D Structural Data Processing: 3221

- 3222 – *Geometric Encoding:* Utilizing encoders 3223
3224 like PointNet, Point Transformer, or Graph 3225
3226 Neural Networks (GNNs) to process irregu- 3227
3228 lar 3D data such as point clouds (*e.g.*, from 3229
3230 LiDAR scans of geological sites) or molec- 3231
3232 ular meshes. 3232
- 3233 – *Incorporating Physical and Chemical Con- 3234
3235 straints:* Enhancing structural understand- 3236
3237 ing by integrating domain-specific knowl- 3238
3239 edge, such as bond lengths and angles in 3240
3241 molecular modeling or material properties 3242
3243 in engineering design, directly into the en- 3244
3245 coding or reasoning process. This ensures 3246
3247 that interpretations are scientifically plausi- 3248
3249 ble. 3250

3251 • Sensor and Time-Series Data: 3252

- 3253 – *Specialized Encoders:* Using Recurrent 3254
3255 Neural Networks (RNNs), 1D CNNs, or 3256
3257 Transformers adapted for time-series data 3258
3259 from various sensors (*e.g.*, temperature, 3259
3260 pressure, EEG, spectroscopy). 3260
3261 – *Temporal Alignment:* Developing meth- 3261
3262 ods to align asynchronous and heteroge- 3262
3263 neous sensor streams with other modalities 3263
3264 like video or textual logs of experimental 3264
3265 events. 3265

3266 H.3 Advanced Cross-Modal Interaction and 3267 3268 Scientific Reasoning 3268

3269 The true power of incorporating more modalities 3269
3270 lies in enabling more sophisticated cross-modal in- 3270
3271 teractions and reasoning capabilities: 3271

3272 • Unified Multimodal Alignment: 3273

- 3274 – *Shared Embedding Space:* Designing a 3274
3275 common embedding space, potentially 3275
3276 through CLIP-style contrastive learning ex- 3276
3277 tended to audio, video, and 3D data, where 3277
3278 features from diverse modalities represent- 3278
3279 ing the same scientific concept are brought 3279
3280 closer together. 3280
3281 – *Dynamic Modality Weighting:* Implement- 3281
3282 ing mechanisms that allow the model to dy- 3282
3283 namically assign importance (weights) to 3283
3284 3284

different modalities based on the specific scientific task, context, or data quality, ensuring robustness and flexibility.

- **Enhanced Scientific Reasoning Modules:**

- *Multimodal Causal Graph Generation:* Constructing more comprehensive causal relationship graphs by leveraging evidence from diverse inputs. For example, observing a video of an experimental setup, listening to audio cues indicating a process stage, and reading instrument outputs (text/tabular) can collectively inform the nodes and edges of a causal graph explaining an observed phenomenon (e.g., "increased gas pressure [sensor] caused by heating element activation [video/log] led to audible valve release [audio]").
- *Generative Reasoning with Richer Evidence:* Empowering language models to generate more nuanced and well-supported scientific explanations, hypotheses, or experimental designs by integrating and citing evidence from audio logs, video segments, 3D visualizations, and sensor data alongside textual knowledge.

H.4 Task-Driven Training Strategies and Future Directions

Realizing the potential of MLLMs with expanded modalities requires innovative training strategies and a forward-looking research agenda:

- **Multi-Task Learning for Comprehensive Understanding:**

- Training models on a combination of tasks, including fundamental modality understanding (e.g., audio event classification, video action recognition, 3D object recognition) alongside higher-level scientific reasoning tasks (e.g., predicting experimental outcomes, inferring material properties from multimodal sensor data).

- **Leveraging Domain-Specific Knowledge Graphs:**

- Incorporating expert-annotated scientific knowledge graphs (ontologies, relational databases) as prior constraints or auxiliary information during training. This can guide the model in interpreting complex multimodal inputs within the correct scientific

framework (e.g., linking a 3D protein structure to its known functions and interaction pathways).

- **Addressing Data Scarcity:**

- Developing techniques for data augmentation, synthetic data generation (e.g., simulating experimental videos or 3D molecular interactions), and leveraging weakly supervised or self-supervised learning methods, given that large-scale, richly annotated multimodal scientific datasets are often scarce.

- **Future Outlook:** The path forward involves creating benchmark datasets that encompass these richer modalities within scientific contexts. Furthermore, research should explore how these expanded sensory capabilities can enable truly interactive MLLMs that can actively participate in the scientific process, for instance, by suggesting modifications to an experimental setup based on real-time video and sensor feedback. Such advancements will be pivotal in unlocking new frontiers in AI-assisted scientific discovery.

I Discussion of Alternative Views

While this paper champions the significant potential of MLLMs to advance scientific reasoning, it is crucial to engage with pertinent alternative viewpoints and concerns that challenge this optimism. Acknowledging these perspectives allows for a more nuanced understanding of the MLLM adoption pathway in science.

I.1 Domain-Specific Models as a Superior Alternative

One argument suggests that *highly specialized, domain-specific models tailored to individual scientific disciplines may outperform general-purpose MLLMs in reasoning tasks*, as indicated in Figure 9 (a). Proponents of this view highlight that scientific reasoning often requires deep domain expertise, nuanced understanding, and customized data processing pipelines that are difficult to replicate in generalized multimodal architectures (Zhang et al., 2024d; Barman et al., 2025). For example, domain-specific models like AlphaFold for protein structure prediction (Jumper et al., 2021) or symbolic computation tools for solving mathematical problems (Mirzadeh et al., 2024; Lu et al., 2021a) excel precisely because of their narrow focus.

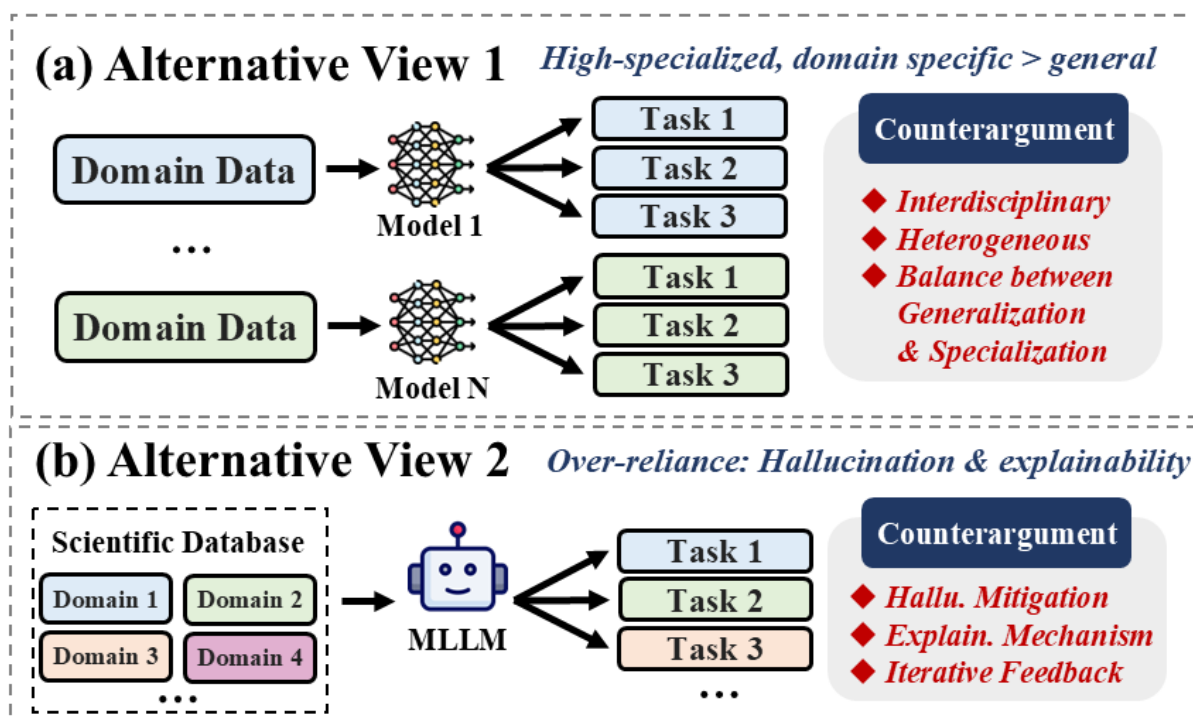


Figure 9: Illustrations of alternative view 1 (a) and view 2 (b), as well as our corresponding counterarguments.

Counterargument. While domain-specific models have demonstrated remarkable success in narrow applications, they lack the flexibility to generalize across multiple domains or integrate diverse modalities of data. Scientific reasoning increasingly involves interdisciplinary approaches, such as bioinformatics combining biology and computational techniques, or climate science requiring both geospatial and textual analysis (Reddy and Shojaee, 2024). MLLMs, by design, offer the ability to process and reason over heterogeneous data sources, enabling a broader and more integrative approach (Zhang et al., 2024f). Furthermore, domain-specific knowledge can still be finetuned within MLLMs, allowing these systems to leverage the best of both generalization and specialization (Chen et al., 2023c).

More Elaboration can be seen as follows:

• **Necessity of Customized Architectures and Data Handling:**

- *What creates the necessity of a customized model?* Many scientific domains rely on unique, highly structured data formats (e.g., SMILES strings in chemistry, FASTA sequences in genomics, PDB files for protein structures, specific notations in physics diagrams or mathematical proofs). General MLLMs, even with multimodal capa-

bilities, may lack the inherent inductive biases or fine-grained understanding to parse, interpret, and reason over these specialized representations with the required precision. A custom model can be built from the ground up to "speak" the native language of that data.

- Domain-specific models can embed established physical laws, chemical reaction rules, or biological pathways more explicitly, either through architectural design or by training on curated datasets that heavily emphasize these principles. This contrasts with MLLMs that must infer these rules implicitly from broader, potentially noisier data.
- The sheer precision required for tasks like predicting quantum mechanical properties, drug-target interaction affinities, or complex engineering tolerances might necessitate models whose entire objective function and training data are geared towards minimizing error in that specific, narrow domain, rather than general plausibility.

• **Knowledge Depth and Nuance:**

- Scientific breakthroughs often require profound depth of knowledge and nuanced understanding within a very specific sub-field.

3411 While MLLMs aim for breadth, they might
3412 only achieve a superficial understanding in
3413 areas requiring years of human expert training.
3414 The "curse of dimensionality" applies
3415 to knowledge as well; covering all scientific
3416 domains deeply in one model is an immense
3417 challenge.

- 3418 – Training datasets for general MLLMs, even
3419 if including scientific papers, might be diluted
3420 by vast amounts of non-scientific text,
3421 or may not capture the long tail of niche,
3422 but critical, scientific knowledge that specialized
3423 databases and models incorporate.

3424 While the counterargument in the main paper highlights
3425 MLLM flexibility for interdisciplinary tasks, proponents
3426 of specialization argue that many critical scientific
3427 problems are bottlenecked by depth in a single
3428 discipline, not breadth. They contend that fine-tuning
3429 an MLLM might offer some domain adaptation, but it's
3430 unlikely to match a model purpose-built for the intricacies
3431 of that domain.

3432 I.2 Risks of Over-reliance on MLLMs

3433 Another valid concern is *the potential over-reliance*
3434 *on MLLMs, which could exacerbate issues such as*
3435 *hallucination and lack of explainability*, as shown in
3436 Figure 9 (b). Critics argue that MLLMs, while powerful,
3437 are prone to generating plausible-sounding but incorrect
3438 or unsubstantiated outputs. This risk is particularly
3439 concerning in scientific reasoning, where accuracy and
3440 rigor are paramount (Bai et al., 2024b). Additionally,
3441 the black-box nature of these models makes it difficult
3442 for researchers to validate their reasoning processes
3443 or trust their conclusions, which could hinder their
3444 adoption in critical scientific applications (Cambria et al.,
3445 2024; Rodis et al., 2024; Dang et al., 2024b).

3446 **Counterargument.** These concerns are valid and
3447 underscore the need for a cautious and measured
3448 approach to MLLM adoption. However, rather than
3449 dismissing MLLMs outright, these issues highlight
3450 areas for improvement. For example, integrating
3451 explainability mechanisms, such as visual attention
3452 maps (Chefer et al., 2021; Dehimi and Tolba, 2024)
3453 or rationale generation (Hu and Yu, 2024; Wu et al.,
3454 2024a), can enhance transparency. Additionally,
3455 hybrid models that combine MLLMs with symbolic
3456 reasoning (Li et al., 2024a; Zhou et al., 2024a) or
3457 expert systems (Guan et al., 2024; Niu et al., 2024)
3458 can mitigate

3461 risks while maintaining strengths of multimodal
3462 reasoning. Finally, iterative feedback loops and
3463 human-in-the-loop systems can ensure reliability
3464 in reasoning workflows (Xiao et al., 2024; Zou
3465 et al., 2024; Zheng et al., 2024).

3466 **More Elaboration** can be seen as follows:

3467 • The Amplified Peril of Scientific Hallucination: 3468

- 3469 – Hallucinations in a general chatbot might
3470 be an annoyance; in science, they can be
3471 dangerous or resource-wasting. This isn't
3472 just about generating factually incorrect
3473 statements, but also subtle misinterpretations
3474 of complex multimodal inputs (*e.g.*,
3475 misreading a crucial detail in a micrograph
3476 or circuit diagram) that lead to flawed
3477 reasoning chains.
- 3478 – The "plausible-sounding but incorrect" nature
3479 of MLLM outputs is especially insidious
3480 in science, where a subtly flawed equation,
3481 a non-existent chemical reaction pathway,
3482 or an incorrectly inferred biological interaction
3483 could derail research or lead to unsafe
3484 experimental protocols.

3485 • The Imperative of Explainability and Verifiability: 3486

- 3487 – The scientific method fundamentally relies
3488 on transparency, reproducibility, and the
3489 ability to scrutinize the reasoning process.
3490 The "black-box" nature of many MLLMs
3491 directly conflicts with this. How does one
3492 verify the multi-step reasoning an MLLM
3493 used to arrive at a novel hypothesis from a
3494 complex set of inputs (*e.g.*, a research paper,
3495 experimental data, and a diagram)?
- 3496 – Current explainability techniques (*e.g.*,
3497 attention maps, rationale generation) often
3498 provide superficial or post-hoc justifications
3499 that may not reflect the true internal
3500 "reasoning" of the model, if such a
3501 coherent process even exists in a human-
3502 understandable form. This makes debugging
3503 errors or building trust in novel, MLLM-
3504 generated insights exceptionally difficult.
3505

3506 • Bias Propagation and Deskilling Concerns:

- 3507 – MLLMs trained on existing scientific literature
3508 and datasets will inevitably in-

3509 herit and potentially amplify existing bi-
3510 ases within those sources (*e.g.*, overrepresent-
3511 ation of certain research topics, demo-
3512 graphic groups in clinical data, or estab-
3513 lished theories at the expense of emerging
3514 ones).

- 3515 – A more philosophical concern is the poten-
3516 tial for "deskilling." If researchers become
3517 overly reliant on MLLMs for hypothesis
3518 generation, data interpretation, or even ex-
3519 perimental design, there's a risk of atrophy-
3520 ing fundamental human scientific reason-
3521 ing and critical thinking skills, reducing sci-
3522 entists to sophisticated prompt engineers.

3523 The main paper's counterargument focuses on mit-
3524 igation strategies. However, critics would ar-
3525 gue that the current state of these mitigations
3526 is insufficient for high-stakes scientific applica-
3527 tions, and the burden of proof lies in demon-
3528 strating consistent reliability, deep interpretability,
3529 and bias control before widespread adoption in
3530 critical scientific endeavors. The very nature of
3531 MLLMs—learning statistical patterns rather than
3532 causal mechanisms—might pose a fundamental
3533 limit to their trustworthiness in a domain that
3534 prizes causal understanding.

3535 **J Clarification of LLM Usage**

3536 In the preparation of this manuscript, the authors
3537 leveraged Gemini-Pro-2.5 as an assistive writing
3538 tool. The use of the LLM was strictly confined
3539 to an editorial capacity, primarily for refining lan-
3540 guage, improving grammatical correctness, and
3541 enhancing the clarity and readability of the text.
3542 The core intellectual contributions including the
3543 central thesis, the proposed four-stage research
3544 roadmap, the analysis of challenges, and the over-
3545 all structure of the paper were conceived and devel-
3546 oped exclusively by the human authors. The au-
3547 thors meticulously reviewed, edited, and approved
3548 all final wording to ensure it accurately reflects
3549 their original ideas and arguments, and bear full
3550 responsibility for the accuracy and integrity of the
3551 content.
3552

Table 6: Summary of scientific (M)LLMs.

Paper	Organization	Venue	Multimodal	Parameter	Domain(s)			
					Math	Physics	Chemistry	Biology
General-Reasoner (Ma et al., 2025)	Waterloo	NeurIPS'25		4B/7B/14B	✓	✓	✓	✓
AI-Reasoner (Tang et al., 2025)	HKU	NeurIPS'25		-	✓	✓	✓	✓
OmniScience (Prabhakar et al., 2025)	SES AI			70B	✓	✓	✓	✓
Galactica (Taylor et al., 2022)	Meta		✓	125M/1.3B/6.7B/30B/120B	✓	✓	✓	✓
LLM-SR (Shojaee et al., 2024)	Virginia Tech			8x7B	✓	✓		✓
SciReasoner (Wang et al., 2025b)	Shanghai AI Lab			1.7B/8B		✓	✓	✓
SciLitLLM (Li et al., 2024k)	USTC			7B/14B		✓	✓	✓
Darwin series (Xie et al., 2023)	University of New South Wales			7B		✓	✓	
SPMM (Chang and Ye, 2022)	KAIST	Nature Communications'24	✓	-			✓	✓
InstructMol (Cao et al., 2023b)	IDEA & HKUST	COLING'25	✓	7B			✓	✓
BioinspiredLLM (Luu and Buehler, 2023)	MIT	Advanced Science'24		13B			✓	✓
nach0 (Livne et al., 2023)	NVIDIA	Chemical Science'24	✓	250M/780M			✓	✓
Mole-BERT (Xia et al., 2023)	Westlake University	ICLR'23		-			✓	✓
BioReason (Fallahpour et al., 2025)	University of Toronto			1B/4B				✓
BioGPT (Luo et al., 2022)	Microsoft	Briefings in Bioinformatics'22		355M				✓
Evolla (Zhou et al., 2025)	Westlake University		✓	80B				✓
Prollama (Lv et al., 2024)	Peking University			7B				✓
ether0 (Narayanan et al., 2025)	FutureHouse Inc	NeurIPS'25		24B			✓	
MOOSE-Chem2 (Yang et al., 2025)	NTU & Shanghai AI Lab	NeurIPS'25		-			✓	
Perovskite-LLM (Liu et al., 2025c)	HKUST(GZ)	EMNLP Findings'25		8B			✓	
Chemdm (Zhao et al., 2024c)	Shanghai Jiaotong University			13B			✓	
ChemDFM-X (Zhao et al., 2024b)	Shanghai Jiaotong University		✓	8B			✓	
ChemLLM (Zhang et al., 2024a)	Shanghai AI Lab			7B			✓	
GIT-Mol (Liu et al., 2023c)	Peng Cheng Lab	Computers in Biology and Medicine'24	✓	700M			✓	
MolGPT (Bagal et al., 2021)	International Institute of Information Technology	Journal of Chemical Information and Modeling'21		6M			✓	
MOOSE-Chem (Yang et al., 2024e)	NTU & Shanghai AI Lab			-			✓	
BatGPT-Chem (Yang et al., 2024c)	Shanghai Jiaotong University			15B			✓	
DARWIN 1.5 (Xie et al., 2024)	University of New South Wales			7B			✓	
MolMetaLM (Wu et al., 2024b)	Central South University			-			✓	
SMTED (Soares et al., 2024)	IBM			289M			✓	
MathCoder2 (Lu et al., 2024b)	CUHK	ICLR'25		7B	✓			
MathCoder-VL (Wang et al., 2025a)	CUHK	ACL'25 Findings	✓	2B/8B	✓			
MathCoder (Wang et al., 2023c)	CUHK	ICLR'24		7B/13B	✓			
MAMmoTH1 (Yue et al., 2023)	UWaterloo	ICLR24		7B/13B/70B	✓			
Math-LLaVA (Shi et al., 2024b)	NUS	EMNLP Findings'24		13B	✓			
JiuZhang 2.0 (Zhao et al., 2023b)	RUC & iFLYTEK	KDD'23		-	✓			
JiuZhang 1.0 (Zhao et al., 2022)	RUC & iFLYTEK	KDD'22		145M	✓			
Minerva (Lewkowycz et al., 2022)	Google	NeurIPS'22		8B/62B/540B	✓			
Hypertree Proof Search (Lample et al., 2022)	Meta	NeurIPS'22		-	✓			
Qwen2.5-Math (Yang et al., 2024b)	Alibaba			1.5B/7B/72B	✓			
Qwen2-Math (Yang et al., 2024a)	Alibaba			1.5B/7B/72B	✓			
Qwen2-Math-Instruct (Yang et al., 2024a)	Alibaba			1.5B/7B/72B	✓			
MathGPT (TAL Education, 2023)	TAL Group		✓	130B	✓			
math-specialized Gemini 1.5 Pro (Team et al., 2024b)	Google		✓	-	✓			
InternLM2-Math (Ying et al., 2024)	Shanghai AI Lab			1.8B/7B/20B/8x22B	✓			
InternLM2.5-StepProver (Wu et al., 2024c)	Shanghai AI Lab			7B	✓			
Llemma (Azerbayev et al., 2023)	Princeton University & Eleuther AI			7B/34B	✓			
ChatGLM-Math (Xu et al., 2024b)	Zhipu AI			32B	✓			
MetaMath (Yu et al., 2023)	Cambridge & Huawei			7B/13B/70B	✓			
MathGLM (Yang et al., 2023b)	Tsinghua & Zhipu AI			10M/100M/335M/500M/2B/6B/10B	✓			
MathGLM-Vision (Yang et al., 2024d)	Tsinghua & Zhipu AI		✓	9B/19B/32B	✓			
Skywork-13B-Math (Yang et al., 2023a)	SkyworkAI		✓	7B/13B	✓			
DeepSeekMath-V2 (Shao et al., 2025)	DeepSeek AI			685B	✓			
DeepSeekMath (Shao et al., 2024)	DeepSeek AI			7B	✓			
DeepSeekProver-V1 (Xin et al., 2024a)	DeepSeek AI			7B	✓			
DeepSeek-Prover-V1.5 (Xin et al., 2024b)	DeepSeek AI			7B	✓			
Mathstral (Mistral AI Team, 2024)	Mistral AI			7B	✓			
JiuZhang 3.0 (Zhou et al., 2024b)	RUC & iFLYTEK			7B/8B	✓			
Math-LLM (Liu et al., 2024c)	East China Normal University		✓	8.26B/7B/72B	✓			
GPT-f (Polu and Sutskever, 2020)	OpenAI			160M/400M/700M/	✓			
Rho-Math (Lin et al., 2024)	Microsoft			1B/7B	✓			
WizardMath (Luo et al., 2023)	Microsoft			7B/70B	✓			
Xwin-LM (Ni et al., 2024)	Microsoft			7B/13B/70B	✓			
MAMmoTH2 (Yue et al., 2024c)	UWaterloo			7B/8B	✓			
GAIRMath-Abel (Chem et al., 2023)	Shanghai Jiaotong University			7B/13B/70B	✓			
KwaiYiiMath (Fu et al., 2023)	Kuaishou			13B	✓			
k0-math (MoonshotAI, 2024)	Moonshot AI			-	✓			
NuminaMath (Beeching et al., 2024)	Numina			7B/72B	✓			