

LLM Agents for Education: Advances and Applications

Anonymous ACL submission

Abstract

Large Language Model (LLM) agents are transforming education by automating complex pedagogical tasks and enhancing both teaching and learning processes. In this survey, we present a systematic review of recent advances in applying LLM agents to address key challenges in educational settings, such as feedback comment generation, curriculum design, etc. We analyze the technologies enabling these agents, including representative datasets, benchmarks, and algorithmic frameworks. Additionally, we highlight key challenges in deploying LLM agents in educational settings, including ethical issues, hallucination and overreliance, and integration with existing educational ecosystems. Beyond the core technical focus, we include in Appendix A a comprehensive overview of domain-specific educational agents, covering areas such as science learning, language learning, and professional development.

1 Introduction

Artificial intelligence techniques are increasingly used in education to enable personalized learning and intelligent tutoring (Chen et al., 2020; Zhai et al., 2021; Pedro et al., 2019). While traditional educational data mining approaches (Shafiq et al., 2022; Khan and Ghosh, 2021; Abdelrahman et al., 2023; Song et al., 2022; Wang et al., 2022; Gao et al., 2021), such as knowledge tracing and cognitive diagnosis, have made significant progress in reshaping the human-learning paradigm by analyzing student behaviors and assessing knowledge states, they still face major challenges in real-world applications. These challenges include shallow contextual understanding, limited interactive capabilities, and difficulties in generating adaptive, personalized learning materials, etc (Zhu et al., 2024; Laak and Aru, 2024; Tan et al., 2023).

The strong natural language understanding of Large Language Models (LLMs) and the task automation capabilities of LLM agents make them

valuable for addressing challenges in education (Weng, 2023; Wang et al., 2024a). First, memory enables LLM agents to retain both long-term knowledge about students’ study habits and short-term context from real-time interactions, enhancing contextual understanding and ensuring personalized learning experiences across various educational tasks (Zhang et al., 2024d). Second, tool use allows LLM agents to access external resources, perform complex calculations, and retrieve real-time information, enabling them to automate intricate educational tasks such as grading, knowledge retrieval, and adaptive content generation, thereby overcoming limited interactivity and enhancing engagement (Gao et al., 2023). Third, planning supports structured learning by decomposing complex topics, predicting optimal learning paths, and dynamically adjusting instructional strategies, allowing LLM agents to autonomously guide students through personalized learning experiences (Huang et al., 2024b).

In addition to these core architectural capabilities, we identify personalization, explainability, and multi-agent communication as essential for effective educational LLM agents, as they enable nuanced instructional behavior and collaborative reasoning. Personalization allows agents to tailor instruction to individual learners’ needs and preferences (Razafinirina et al., 2024). Explainability enables them to provide interpretable justifications for feedback and decisions (Abu-Rasheed et al., 2024b). Multi-agent communication facilitates role-based collaboration, such as between a planner and a critic, to improve robustness and task coverage (Wu et al., 2023a). By integrating these features, LLM agents enhance understanding and engagement while streamlining educational workflows for greater adaptability and efficiency.

In this survey, we provide a comprehensive review of LLM agents in educational settings, with a focus on their underlying technical foundations

and general-purpose pedagogical capabilities. We begin by introducing the core abilities of educational LLM agents, including memory, tool use, planning, personalization, explainability, and multi-agent communication, and discuss their potential to automate and enhance diverse educational tasks. Given the highly applied nature of the education domain, we propose a task-centric taxonomy in Figure 1 that organizes recent advances around core educational tasks. We categorize educational LLM agents based on their roles in supporting teachers and students, capturing both instructional and learning-focused functions: (1) Teaching Assistance Agents, which support educators by automating tasks such as *Classroom Simulation (CS)*, *Feedback Comment Generation (FCG)*, and *Curriculum Design (CD)*; and (2) Student Support Agents, which facilitate personalized learning through *Adaptive Learning (AL)*, *Knowledge Tracing (KT)*, and *Error Correction and Detection (ECD)*. The overview of LLM agents for education and illustrative examples of each task are presented in Figure 2. We further highlight critical challenges in deploying these agents, including ethical issues, hallucination and overreliance, and integration into existing educational ecosystems. Finally, we compile essential datasets, benchmarks, and evaluation protocols to support future research on LLM agent-based educational systems. We summarize our contributions as follows:

- **Novel task-centric taxonomy.** We propose a structured taxonomy that categorizes LLM agents based on core educational tasks, highlighting roles in teaching assistance and student support to unify analysis across applications.
- **Current challenges and future directions.** We analyze critical challenges that need to be addressed for the effective deployment of LLM agents for education, including issues related to ethical issues, hallucination, and integration into real-world educational ecosystems.
- **Compilation of essential resources.** We compile comprehensive datasets and benchmarks to support future research efforts and facilitate the development of more robust and effective LLM-driven educational solutions.¹

Beyond the main technical focus, Appendix A provides an overview of domain-specific educa-

tional agents, including applications in science learning, language learning, and professional development. We outline domain-specific challenges and review recent advances, benchmarks, and datasets. Readers are encouraged to refer to the appendix for further details.

2 LLM Agents for Education

LLM agents have a set of connected abilities that help them handle complex tasks and provide meaningful support in education. These core features work together to let LLM agents do more than just find information—they can interact in ways that are flexible, responsive, and tailored to each learner. Their main strengths in education include strong memory, the ability to use tools, planning skills, personalization, clear explanations, and the ability to work with other agents. Specifically,

Memory. Memory in LLM agents includes long-term (foundational knowledge, e.g., commonsense) and short-term (current interaction data) components (Sumers et al.). Active management via summarization and retrieval ensures relevant context is maintained (Chen et al., 2023; Liang et al., 2023; Zhong et al., 2023). This allows agents to track student progress and personalize responses, though sophisticated filtering is crucial to maintain the quality against noisy interaction data.

Tool Use. To overcome limitations like knowledge cutoffs or calculation difficulties, LLM agents use external tools such as search engines, databases, or APIs (Qin et al., 2023; Zhuang et al., 2023). This expands their functionality, providing access to current information and specialized capabilities. By integrating these tools, LLM agents expand their functional capabilities, ensuring that they can provide accurate and relevant information while supporting diverse educational tasks.

Planning. Planning allows agents to actively support learning by breaking down complex goals into smaller, manageable steps and adjusting based on how the student is doing (Yao et al., 2023; Valmeekam et al., 2023). This means understanding the learning goals, creating step-by-step plans, personalizing the learning path, and making changes based on feedback and progress. For longer-term tasks, LLMs can act as “meta-controllers,” using methods like “Pedagogical Steering” to stay aligned with teaching goals, which thus leads to dynamic, emergent curricula co-created with students (Zhang et al., 2025c).

¹Due to the page limit, we present more details in Appendix B

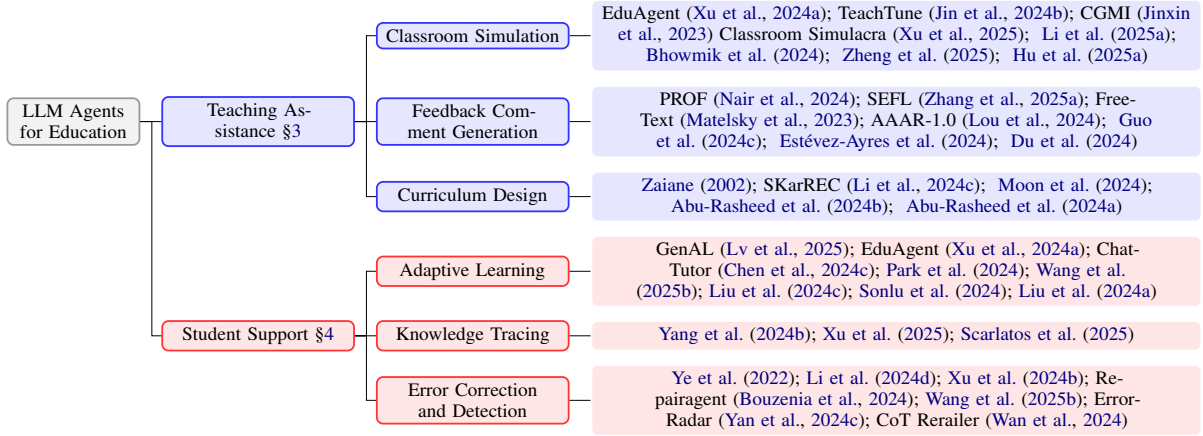


Figure 1: Taxonomy of representative research on LLM agents for education.

Personalization. Personalization is a hallmark of effective education, and LLM agents excel in this area by adapting to individual learning styles and needs (Chen et al.). They can serve as tutors, teaching assistants, or even simulate peer interactions to enhance the learning experience (Jin et al., 2024b; Guo et al., 2024c). Additionally, LLM agents can simulate user behaviors in recommendation systems, which can be applied to personalize educational content (Li et al., 2024c; Chu et al., 2025). Such a technology of personalization transforms traditional education, making it more accessible and tailored to each student’s unique needs.

Explainability. Explainability is crucial for building trust and facilitating learning in educational settings (Wu et al., 2024b). LLM agents must provide clear, step-by-step explanations that are understandable to students. Clear explanations are particularly important in subjects like STEM, where step-by-step reasoning is essential for student understanding (Nair et al., 2024; Zhang et al., 2025a). Consequently, the focused effort required to make educational LLM agents more transparent and their decision-making processes more scrutable.

Multi-Agent Communication. LLM agents can facilitate richer interactions in collaborative learning environments (Chen et al.; Wu et al., 2023a). Multiple agents could coordinate on group projects or simulate peer learning, fostering critical thinking and diverse perspectives (Jinxin et al., 2023; Yue et al., 2024). This approach could support complex simulations of teamwork and real-world professional settings, especially when guided by new teaching frameworks for organization and evaluation (Xu et al., 2024a).

Table 1 summarizes the core capabilities re-

quired by each educational task supported by LLM agents. For each task, we highlight the primary capabilities identified from the reviewed literature. However, certain studies have also explored additional capabilities beyond these primary ones to enhance specific tasks—for example, Guo et al. (2024c) leverage multi-agent communication to improve feedback correctness. We omit these secondary capabilities from the table to more clearly illustrate the primary focus areas associated with each task.

3 Agent for Teaching Assistance

Agents for teaching assistance are designed to support educators in the learning environment. These agents leverage LLMs to provide personalized, scalable, and efficient support across various aspects of the educational process. Their primary objectives are to enhance teaching quality, enrich student learning experiences, and reduce educators’ workload. By incorporating advanced capabilities, LLM agents can effectively support key educational tasks, including classroom simulation (§3.1), feedback comment generation (§3.2), and curriculum design (§3.3).

3.1 Classroom Simulation

Classroom simulation refers to the ability of teaching agents to replicate and model various classroom scenarios, such as student-teacher dialogues, collaborative learning activities, and problem-solving tasks. These simulations create dynamic and interactive learning environments where educators can experiment with different teaching strategies, assess student reactions, and receive real-time feedback on how various pedagogical approaches may

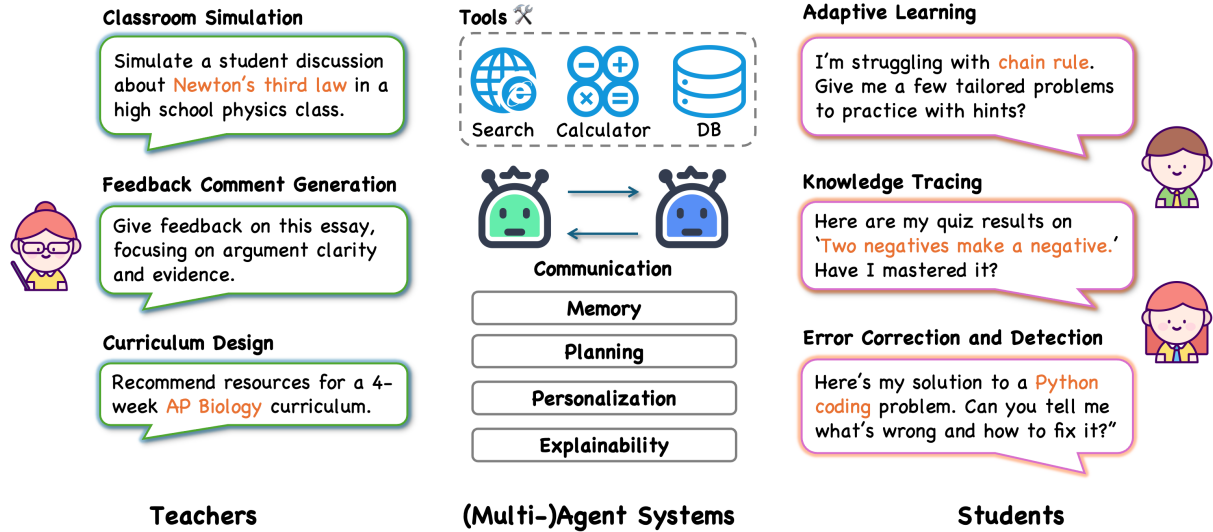


Figure 2: The overview of LLM Agents for education. Teachers and students interact with LLM agents by submitting task-specific prompts. The agents respond using core capabilities such as memory, planning, tool use, and personalization to carry out tasks that support instruction and learning.

| Task | Memory | Tool Use | Planning | Personalization | Explainability | Multi-Agent Comm. |
|------|--------|----------|----------|-----------------|----------------|-------------------|
| CS | ✓ | — | ✓ | — | — | ✓ |
| FCG | — | — | — | ✓ | ✓ | — |
| CD | ✓ | ✓ | ✓ | ✓ | — | — |
| AL | ✓ | — | ✓ | ✓ | — | — |
| KT | ✓ | — | — | ✓ | — | — |
| ECD | — | ✓ | — | ✓ | ✓ | — |

Table 1: Mapping of educational tasks to essential LLM agent capabilities. ✓ indicates a primary capability; others may also be used as supporting components in some systems.

unfold. By simulating these classroom dynamics, educators can refine their methods, anticipate student challenges, and enhance overall instructional effectiveness, all without the constraints of a physical classroom setting.

As shown in Table 1, effective classroom simulation relies on memory, planning, and multi-agent communication to accurately model student behavior. Previous studies (Xu et al., 2024a; Jin et al., 2024b; Li et al., 2025a) demonstrate that LLM-based agents can predict fine-grained student behaviors across diverse personas and past learning patterns, aligning closely with real teachers' expectations. To enhance simulation, the CGMI framework (Jinxin et al., 2023) uses a tree-based cognitive architecture with memory, reflection, and planning modules to simulate roles like teacher, student, and supervisor, improving realism. Similarly, Classroom Simulacra (Xu et al., 2025) incorporates a transferable iterative reflection module for more accurate behavior simulation. These systems enable automated interactions that reduce educators'

task loads while broadening the exploration of student profiles. Simulations can also test educational strategies tailored to different profiles, enhancing teaching quality, as shown in studies by Bhowmik et al. (2024) and Zheng et al. (2025). Additionally, Hu et al. (2025a) demonstrate how LLMs can refine teaching plans through integrated simulations. To sum up, classroom simulations can be leveraged to test different educational strategies tailored to diverse student profiles, ultimately enhancing the quality of education.

3.2 Feedback Comment Generation

Providing timely, relevant, and constructive feedback is a cornerstone of effective education. Teaching assistance agents can generate automated feedback comments on students' assignments, quizzes, and projects.

For example, Guo et al. (2024c) leverage multi-agent communication to provide accurate feedback to students through a two-agent system. Specifically, Agent 1 generates initial feedback based on the students' responses, while Agent 2 eval-

uates and refines this feedback to prevent overpraise and excessive inferences. Furthermore, [Nair et al. \(2024\)](#) design a training strategy called PROF, which trains an automated LLM-based writing comment generator through reinforcement learning. This system adopts an iterative pipeline to simulate various student writing styles and incorporates a more advanced revision model (e.g., GPT-4) to provide the quality of the feedback as rewards. Similarly, SEFL ([Zhang et al., 2025a](#)) enhances feedback generation by having LLM agents role-play both students and teachers to generate data, which is then used to fine-tune models and improve feedback capabilities. These systems have also been deployed in real-world applications, such as FreeText ([Matelsky et al., 2023](#)), which pairs student responses with teacher-provided criteria, enabling the agent to identify strengths and weaknesses and provide targeted feedback for improvement. Beyond traditional feedback, advanced LLM agents are now capable of handling more complex, expertise-intensive tasks. For example, [Du et al. \(2024\)](#) explores the potential of LLM agents as assistants for natural language processing paper reviewing tasks, while AAAR-1.0 ([Lou et al., 2024](#)) evaluates agents' capabilities in areas such as equation inference, experiment design, paper weakness analysis, and review critique, revealing their potential in conducting research tasks.

However, another line of research highlights that agent-generated feedback still faces challenges in handling complex tasks, such as programming and the review of professional academic papers ([Estévez-Ayres et al., 2024](#); [Lou et al., 2024](#)). For instance, these agents may struggle with concepts like starvation and deadlocks, leading to inaccurate or incorrect feedback. Future work could focus on integrating external tools (e.g., search engines) and enhancing memory mechanisms to better address complex problem-solving scenarios, as well as refining the personalization of feedback across diverse learning contexts.

3.3 Curriculum Design

To ensure that students follow personalized learning paths aligned with their knowledge level and domain, it is crucial to develop effective curriculum design strategies. This is a complex task that relies on multiple LLM agent capabilities, including memory, tool use, planning, personalization, and explainability.

[Zaiane \(2002\)](#) first introduces recommendation systems into e-learning, utilizing web mining techniques to suggest online learning activities or shortcuts on course websites. The integration of LLM-based agents has recently enabled more sophisticated curriculum design through dynamic sequencing and content adaptation. Curricula can be built using either retrieval-based or generation-based methods. Retrieval-based methods involve agents accessing a database or their own memory to suggest existing resources—such as textbooks, research papers, or online content—based on student queries, past behavior, or content similarities ([Shahzad et al., 2025](#); [Li et al., 2024c](#)). In contrast, generative methods create new learning content tailored to an individual student's learning style, knowledge gaps, and interests ([Moon et al., 2024](#)). Moreover, to enhance the understanding and acceptance of recommended content, these agents need to provide reasons for their recommendations. Explaining the rationale behind suggestions fosters trust and enables students to make more informed decisions about the resources they engage with. For example, [Abu-Rasheed et al. \(2024a,b\)](#) incorporate knowledge graphs to guide curriculum pathways toward curated and trustworthy sources. This approach not only enhances the interpretability of recommendations but also reduces the risk of generating misinformation, thereby improving the quality of the learning experience.

Looking ahead, future work should focus on integrating agents with adaptive learning systems to dynamically adjust recommendations in response to real-time student performance. In addition, a hybrid approach that combines generative and retrieval-based methods may improve both the accuracy and diversity of curriculum content. Also, incorporating multimodal resources such as interactive media, video, and immersive simulations could further enhance the learning experience.

4 Agent for Student Support

LLM agent-based student support systems aim to provide real-time personalized assistance without requiring direct teacher involvement. Unlike traditional rule-based systems, these agents offer interactive and adaptive feedback, enabling students to progress at their own pace. Core educational tasks in student support include adaptive learning (§4.1), knowledge tracing (§4.2), and error correction and detection (§4.3).

4.1 Adaptive Learning

LLM agents offer the potential to build self-sustaining adaptive learning systems that operate without direct teacher involvement. These systems dynamically tailor instruction based on student performance, enabling personalized learning at scale. Building on the abilities of memory and personalization, agents are able to maintain and update a structured representation of the learner, commonly referred to as a student profile. This profile informs content selection, pacing, and feedback strategies.

Several implementations exemplify this adaptive approach. GenAL (Lv et al., 2025) integrates external tools such as automated programs and web searches to construct comprehensive student profiles and inform instructional planning. Based on these profiles, the agent assigns tasks that align with the student’s current knowledge level and updates its memory dynamically. EduAgent (Xu et al., 2024a) introduces a structured profiling mechanism comprising four distinct cognitive patterns: gaze behavior linked to physiological memory, motor behavior mapped to motor memory, cognitive state associated with cognitive memory, and post-course assessments contributing to knowledge memory. This structured representation enhances adaptive decision-making by providing a multi-faceted view of student learning states. Chen et al. (2024c) propose a system consisting of interaction, reflection, and reaction, with each component composed of specific LLM tools and memory modules. Furthermore, a meta-agent is introduced to control the information flow through these agents.

Some recent research focuses on modeling students through modular memory or state components that capture *cognitive*, *affective*, and *psychological* dimensions. These representations are often formalized within a state–action framework, where the state space encodes learner traits and the action space governs instructional adaptations. A typical cognitive state representation includes tracking a student’s knowledge proficiency, comprehension levels, and misconceptions (Park et al., 2024; Liu et al., 2024c,a), allowing the agent to tailor explanations, adjust difficulty levels, and reinforce concepts dynamically. Recent studies highlight the importance of affective state modeling, as emotional factors such as motivation, interest, and self-efficacy significantly influence learning outcomes. For instance, Park et al. (2024) propose an affective state model that enables agents to ad-

just feedback tone, provide encouragement, and regulate pacing to maintain engagement. Another crucial dimension of adaptation involves learning preferences and personality traits. These studies integrate personality with memory design, tracking the psychological state of the students. Wang et al. (2025b) integrate learning preferences into state modeling, recognizing that students process information differently depending on instructional format and modality. Adapting content to these preferences enhances retention and learning efficiency. Moreover, Liu et al. (2024c) apply the Big Five personality model (Roccas et al., 2002) to personalize tutoring strategies, acknowledging that individual differences shape learning experiences (Sonlu et al., 2024).

Emerging works also explore multi-agent systems for adaptive learning, where specialized agents collaborate to enhance personalization. For example, Wang et al. (2025b) design five agents—Gap Identifier, Learner Profiler, Dynamic Learner Simulator, Learning Path Scheduler, and Content Creator—to deliver goal-oriented, personalized instruction. Similarly, OATutor (Pardos et al., 2023) provides an experimental platform for modular adaptive learning, allowing researchers to design scalable, domain-general tutoring agents.

4.2 Knowledge Tracing

Knowledge tracing is essential for monitoring a learner’s evolving understanding and predicting future performance. While traditional methods use statistical or deep learning models to estimate mastery, LLM agents offer a more dynamic, personalized approach by leveraging natural language understanding and adaptive instruction.

Recent advancements have explored multi-agent frameworks for knowledge tracing. For instance, Yang et al. (Yang et al., 2024b) propose a multi-agent system with three specialized agent roles: administrator, judge, and critic. In this framework, the administrator delegates knowledge tracing tasks to judges, who collaborate through discussions to assess the student’s cognitive state. The critic agent then evaluates the outcome and determines whether the assessment criteria are met, ensuring a structured yet flexible knowledge tracing process. Other agent-based approaches explore alternative strategies for modeling student knowledge. Xu et al. (2024a) propose simulating students as different personas, allowing agents to adaptively trace knowledge progression based on varied learning

profiles. Meanwhile, [Scarlatos et al. \(2025\)](#) employ dialogue-driven interactions to probe students’ conceptual boundaries, using conversational exchanges to refine knowledge estimation dynamically.

4.3 Error Correction and Detection

Error detection and correction help students refine their understanding through real-time, context-aware feedback. LLM agents can identify mistakes across domains such as writing, programming, and math, and adapt feedback to the learner’s proficiency, acting as intelligent reviewers and assistants ([Ye et al., 2022](#); [Li et al., 2024d](#)).

Agent-based systems leverage state representations and adaptive inference mechanisms to track error patterns and misconceptions dynamically ([Park et al., 2024](#); [Liu et al., 2024c](#); [Wang et al., 2025b](#); [Bouzenia et al., 2024](#)). Recent advancements extend this capability into the multi-modal domain, incorporating direct analysis of student-generated drafts. [Xu et al. \(2024b\)](#) propose a multi-modal LLM framework that processes handwritten or digitally drafted student work. The system first extracts and converts draft content into natural language, enabling the agent to interpret and analyze handwritten responses. The agent then provides indirect yet effective instructional feedback, guiding students toward self-correction and deeper comprehension. Moreover, CoT Rerailer ([Wan et al., 2024](#)) designs a derailment identification process and a rerailment process to conduct error detection when solving math questions. [Zhang et al. \(2025d\)](#) propose the MathCCS benchmark and introduce a sequence error analysis framework that leverages multi-agent collaboration. As the first benchmark for multimodal error detection, ErrorRadar ([Yan et al., 2024c](#)) provides a valuable data foundation for developing multimodal agents in this task.

5 Challenges and Future Directions

In this section, we discuss key challenges that must be addressed to ensure the effective, reliable, and ethical deployment of LLM agents in educational settings. We focus on three critical areas: (1) privacy, bias, and fairness; (2) hallucination and overreliance; and (3) integration with existing educational ecosystems. For each of these challenges, we outline potential research directions aimed at improving the robustness, trustworthiness, and practical applicability of LLM agents in real-world education environments.

5.1 Privacy, Bias and Fairness

Analysis. LLM agents process vast datasets, often containing sensitive personal information, leading to potential privacy risks. Studies highlight low technological readiness and insufficient privacy measures in educational contexts ([Yan et al., 2024a](#)). Emerging research ([He et al., 2024a](#); [Gan et al., 2024](#); [Zhang et al., 2024c](#); [Hua et al., 2024](#); [Huo et al., 2025](#); [Chen et al., 2025b](#)) underscores new privacy and security concerns, emphasizing the need for stronger data protection mechanisms ([Huang, 2023](#); [Ismail, 2025](#); [Khan, 2024](#)). Additionally, bias in LLMs remains a pressing concern, as models trained on large datasets can inadvertently reinforce stereotypes and disparities, affecting educational fairness. Recent work calls for bias mitigation strategies to promote equitable learning experiences ([Adewumi et al., 2024](#); [Aird et al., 2024](#); [Mehrotra et al., 2024](#)). Addressing these biases is essential to ensuring inclusive, unbiased educational outcomes.

Directions. To overcome the above issues, a number of future directions can be explored: (i) *Unlearning for privacy preservation*: leverage advances in machine unlearning ([Liu et al., 2025](#)) to enable agents to retain useful knowledge while selectively forgetting sensitive user data when required. (ii) *Bias detection and mitigation*: develop automated fairness-checking models that evaluate real-time content generated by LLM agents to detect biased explanations, language, or examples. (iii) *Culturally adaptive LLM Agents for global education*: train multilingual, culturally aware educational agents that dynamically adjust explanations based on regional learning norms, historical perspectives, and diverse curricula.

5.2 Hallucination and Overreliance

Analysis. The “hallucination” phenomenon, in which LLMs generate plausible but incorrect or nonsensical information, poses a significant challenge to their reliability in educational contexts ([Zhang et al., 2023](#)). Such inaccuracies can mislead learners by presenting false information with confident and authoritative language, making errors difficult to detect and potentially leading to misconceptions ([da Silva et al., 2024](#); [Jho, 2024](#); [Ho et al., 2024](#)). For example, AI-generated content may fabricate historical events or scientific facts that students unknowingly accept as true. This risk is further amplified by overreliance, as students

and educators may accept AI-generated responses without sufficient critical evaluation. Research has shown that excessive dependence on AI systems can hinder skill acquisition and reduce meaningful engagement with learning materials (Milano et al., 2023; Krupp et al., 2024; Adewumi et al., 2023).

Directions. Some directions can be explored to mitigate hallucinations in LLM agents for education: (i) *Self-correcting AI tutors*: develop LLM agents with self-reflection capabilities (Renze and Guven, 2024), where models review, verify, and refine their own generated content before presenting it to students. (ii) *Hybrid Human-AI feedback loops for educational content verification*: develop teacher-in-the-loop AI systems where educators can review and correct AI-generated responses, refining agents performance over time. (iii) *Pedagogical-aware educational agents*: design agentic frameworks aligned with human pedagogical expertise to mitigate overreliance on AI-generated content.

5.3 Integration with Existing Educational Ecosystems

Analysis. Although LLM agents hold great potential for automating educational practices, it is essential to consider how they can be effectively integrated into existing human-centered educational paradigms. One major challenge is the lack of structured frameworks for integrating LLM agents into educational systems. While models like the FOKE framework (Hu and Wang, 2024) combine foundation models, knowledge graphs, and prompt engineering to provide interactive and explainable learning services, broader adoption requires scalable models that can be validated in diverse real-world educational settings. Additionally, LLMs have been explored as tools to enhance creativity and collaboration in project-based learning (PBL), supporting students through brainstorming, problem-solving, and project execution. However, studies indicate that their effectiveness is limited by the absence of structured guidance frameworks that help educators and students seamlessly incorporate LLM agents into PBL workflows (Zha et al., 2024). Another critical challenge is ensuring equitable access to LLM-powered educational tools, particularly in underfunded schools and institutions with limited AI infrastructure. Platforms such as AI-VERDE (Mithun et al., 2025) aim to democratize access by providing a unified *LLM-as-a-platform* service with built-in access control, privacy-preserving mechanisms, and budget man-

agement. However, achieving widespread adoption still depends on scalable and cost-effective deployment strategies that can support educational institutions at different resource levels.

Directions. Future research should focus on developing standardized frameworks to guide the structured deployment of LLM agents in personalized learning, PBL, and assessment. For example, expanding models like FOKE with adaptive learning strategies, multimodal content, and real-time feedback could enhance instructional effectiveness. Additionally, integrating interactive AI tutors that support student collaboration, project tracking, and contextual guidance would further improve PBL applications. To promote equitable access, develop cost-effective AI tutors through cloud-based and decentralized models would make LLM-powered learning tools more accessible to a wider range of institutions. Finally, to support meaningful integration of LLM agents into educational ecosystems, future work should move beyond task accuracy and explore more practical metrics such as learning gains, user trust, and engagement, which calls for the development of novel education-oriented benchmarks and datasets.

6 Conclusion

In this survey, we presented a comprehensive review of LLM agents for education, focusing on their technical foundations and their potential to transform personalized learning, intelligent tutoring, and pedagogical automation. We proposed a task-centric taxonomy that categorizes LLM agents into Teaching Assistance and Student Support, highlighting their core capabilities such as memory augmentation, tool use, planning, and personalization. We also examined key research challenges, including ethical issues, hallucination and overreliance, and integration with existing educational ecosystems, which must be addressed to ensure reliable and ethical deployment. To support continued progress in this field, we compiled critical datasets, benchmarks, and evaluation methodologies. As LLM agents continue to evolve, their influence on education will expand. Realizing their full potential, however, will require both technical rigor and thoughtful system design. We hope this survey provides a solid foundation for future research and the development of AI-powered educational systems.

Limitations

Considering the rapid development of LLM agents for education, it is possible that some of the most recent advancements may not have been captured at the time of writing. Nevertheless, we have made every effort to ensure that all foundational and representative works are included to provide a comprehensive and accurate overview of the field.

References

Mahyar Abbasian, Iman Azimi, Amir M Rahmani, and Ramesh Jain. 2023. Conversational health agents: A personalized llm-powered agent framework. *arXiv preprint arXiv:2310.02374*.

Alaa Abd-Alrazaq, Rawan AlSaad, Dari Alhuwail, Arfan Ahmed, Pdraig Mark Healy, Syed Latifi, Sarah Aziz, Rafat Damseh, Sadam Alabed Alrazak, Javaid Sheikh, and 1 others. 2023. Large language models in medical education: opportunities, challenges, and future directions. *JMIR Medical Education*, 9(1):e48291.

Ghodai Abdelrahman, Qing Wang, and Bernardo Nunes. 2023. Knowledge tracing: A survey. *ACM Computing Surveys*, 55(11):1–37.

Hasan Abu-Rasheed, Mohamad Hussam Abdulsalam, Christian Weber, and Madjid Fathi. 2024a. Supporting student decisions on learning recommendations: An llm-based chatbot with knowledge graph contextualization for conversational explainability and mentoring. *arXiv preprint arXiv:2401.08517*.

Hasan Abu-Rasheed, Christian Weber, and Madjid Fathi. 2024b. Knowledge graphs as context sources for llm-based explanations of learning recommendations. In *2024 IEEE Global Engineering Education Conference (EDUCON)*, pages 1–5. IEEE.

Adebowale Jeremy Adetayo, Mariam Oyinda Aborisade, and Basheer Abiodun Sanni. 2024. Microsoft copilot and anthropic claude ai in education and library service. *Library Hi Tech News*.

Tosin Adewumi, Lama Alkhalel, Claudia Buck, Sergio Hernandez, Saga Brilieth, Mkpe Kekung, Yelvin Ragimov, and Elisa Barney. 2023. Procot: Stimulating critical thinking and writing of students through engagement with large language models (llms). *arXiv preprint arXiv:2312.09801*.

Tosin Adewumi, Lama Alkhalel, Namrata Gurung, Goya van Boven, and Irene Pagliai. 2024. Fairness and bias in multimodal ai: A survey. *arXiv preprint arXiv:2406.19097*.

Amanda Aird, Paresha Farastu, Joshua Sun, Elena Stefancová, Cassidy All, Amy Volda, Nicholas Mattei, and Robin Burke. 2024. Dynamic fairness-aware recommendation through multi-agent social choice.

ACM Transactions on Recommender Systems, 3(2):1–35.

Antonis Antoniadis, Albert Örwall, Kexun Zhang, Yuxi Xie, Anirudh Goyal, and William Wang. 2024. Swe-search: Enhancing software agents with monte carlo tree search and iterative refinement. *arXiv preprint arXiv:2410.20285*.

Junyi Ao, Yuancheng Wang, Xiaohai Tian, Dekun Chen, Jun Zhang, Lu Lu, Yuxuan Wang, Haizhou Li, and Zhizheng Wu. 2024. *Sd-eval: A benchmark dataset for spoken dialogue understanding beyond words*.

Akari Asai, Jacqueline He, Rulin Shao, Weijia Shi, Amanpreet Singh, Joseph Chee Chang, Kyle Lo, Luca Soldaini, Sergey Feldman, Mike D’arcy, and 1 others. 2024. Openscholar: Synthesizing scientific literature with retrieval-augmented lms. *arXiv preprint arXiv:2411.14199*.

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

Raluca Balan, Anca Doborean, and Costina R Poetar. 2024. Use of automated conversational agents in improving young population mental health: a scoping review. *NPJ Digital Medicine*, 7(1):75.

Kristian G Barman, Sascha Caron, Emily Sullivan, Henk W de Regt, Roberto Ruiz de Austri, Mieke Boon, Michael Färber, Stefan Fröse, Faegheh Hasibi, Andreas Ipp, and 1 others. 2025. Large physics models: Towards a collaborative approach with large language models and foundation models. *arXiv preprint arXiv:2501.05382*.

Patrick Bassner, Eduard Frankford, and Stephan Krusche. 2024. Iris: An ai-driven virtual tutor for computer science education. In *Proceedings of the 2024 on Innovation and Technology in Computer Science Education V. 1*, pages 394–400.

Jackylyn Beredo and Ethel Ong. 2021. Beyond the scene: A comparative analysis of two storytelling-based conversational agents. In *Proceedings of the Asian CHI Symposium 2021*, pages 189–195.

Manojit Bhattacharya, Soumen Pal, Srijan Chatterjee, Sang-Soo Lee, and Chiranjib Chakraborty. 2024. Large language model to multimodal large language model: A journey to shape the biological macromolecules to biological sciences and medicine. *Molecular Therapy-Nucleic Acids*, 35(3).

Saptarshi Bhowmik, Luke West, Alex Barrett, Nuodi Zhang, Chih-Pu Dai, Zlatko Sokolij, Sherry Southerland, Xin Yuan, and Fengfeng Ke. 2024. Evaluation of an llm-powered student agent for teacher training. In *European Conference on Technology Enhanced Learning*, pages 68–74. Springer.

- Ben Bogin, Kejuan Yang, Shashank Gupta, Kyle Richardson, Erin Bransom, Peter Clark, Ashish Sabharwal, and Tushar Khot. 2024. Super: Evaluating agents on setting up and executing tasks from research repositories. *arXiv preprint arXiv:2409.07440*.
- Islem Bouzenia, Premkumar Devanbu, and Michael Pradel. 2024. Repairagent: An autonomous, llm-based agent for program repair. *arXiv preprint arXiv:2403.17134*.
- Cameron Brown and Laura Cruz Castro. 2025. Coordinate: A virtual classroom management tool for large computer science courses using discord. In *Proceedings of the 56th ACM Technical Symposium on Computer Science Education V. 1*, pages 165–171.
- Yan Cai, Linlin Wang, Ye Wang, Gerard de Melo, Ya Zhang, Yanfeng Wang, and Liang He. 2024. Med-bench: A large-scale chinese benchmark for evaluating medical large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17709–17717.
- Ruisheng Cao, Fangyu Lei, Haoyuan Wu, Jixuan Chen, Yeqiao Fu, Hongcheng Gao, Xinzhuang Xiong, Hanchong Zhang, Wenjing Hu, Yuchen Mao, and 1 others. 2024. Spider2-v: How far are multimodal agents from automating data science and engineering workflows? *Advances in Neural Information Processing Systems*, 37:107703–107744.
- Justine Cassell. 2022. Socially interactive agents as peers. In *The Handbook on Socially Interactive Agents: 20 years of Research on Embodied Conversational Agents, Intelligent Virtual Agents, and Social Robotics Volume 2: Interactivity, Platforms, Application*, pages 331–366.
- Haw-Shiuan Chang, Hwai-Jung Hsu, Kuan-Ta Chen, and 1 others. 2015. Modeling exercise relationships in e-learning: A unified approach. In *EDM*, pages 532–535.
- Angxuan Chen, Yuang Wei, Huixiao Le, and Yan Zhang. 2024a. Learning-by-teaching with chatgpt: The effect of teachable chatgpt agent on programming education. *arXiv preprint arXiv:2412.15226*.
- Celia Chen and Alex Leitch. 2024. Llms as academic reading companions: Extending hci through synthetic personae. *arXiv preprint arXiv:2403.19506*.
- Guhong Chen, Liyang Fan, Zihan Gong, Nan Xie, Zixuan Li, Ziqiang Liu, Chengming Li, Qiang Qu, Shiwen Ni, and Min Yang. 2024b. Agentcourt: Simulating court with adversarial evolvable lawyer agents. *arXiv preprint arXiv:2408.08089*.
- Howard Chen, Ramakanth Pasunuru, Jason Weston, and Asli Celikyilmaz. 2023. Walking down the memory maze: Beyond context limit through interactive reading. *arXiv preprint arXiv:2310.05029*.
- Huajun Chen, Keyan Ding, Jing Yu, Junjie Huang, Yuchen Yang, and Qiang Zhang. 2025a. Scitoolagent: A knowledge graph-driven scientific agent for multi-tool integration.
- Jiangjie Chen, Xintao Wang, Rui Xu, Siyu Yuan, Yikai Zhang, Wei Shi, Jian Xie, Shuang Li, Ruihan Yang, Tinghui Zhu, and 1 others. From persona to personalization: A survey on role-playing language agents. *Transactions on Machine Learning Research*.
- Junkai Chen, Zhijie Deng, Kening Zheng, Yibo Yan, Shuliang Liu, PeiJun Wu, Peijie Jiang, Jia Liu, and Xuming Hu. 2025b. Safeeraser: Enhancing safety in multimodal large language models through multimodal machine unlearning. *arXiv preprint arXiv:2502.12520*.
- Lijia Chen, Pingping Chen, and Zhijian Lin. 2020. Artificial intelligence in education: A review. *Ieee Access*, 8:75264–75278.
- Yulin Chen, Ning Ding, Hai-Tao Zheng, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. 2024c. Empowering private tutoring by chaining large language models. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management*, pages 354–364.
- Ziru Chen, Shijie Chen, Yuting Ning, Qianheng Zhang, Boshi Wang, Botao Yu, Yifei Li, Zeyi Liao, Chen Wei, Zitong Lu, and 1 others. 2024d. Scienceagentbench: Toward rigorous assessment of language agents for data-driven scientific discovery. *arXiv preprint arXiv:2410.05080*.
- Shanbo Cheng, Zhichao Huang, Tom Ko, Hang Li, Ningxin Peng, Lu Xu, and Qini Zhang. 2024. Towards achieving human parity on end-to-end simultaneous speech translation via llm agent. *arXiv preprint arXiv:2407.21646*.
- Neeraj Cherakara, Finny Varghese, Sheena Shabana, Nivan Nelson, Abhiram Karukayil, Rohith Kulothungan, Mohammed Afil Farhan, Birthe Nessel, Meriam Moujahid, Tanvi Dinkar, Verena Rieser, and Oliver Lemon. 2023. *FurChat: An embodied conversational agent using LLMs, combining open and closed-domain dialogue with facial expressions*. In *Proceedings of the 24th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 588–592, Prague, Czechia. Association for Computational Linguistics.
- Zhendong Chu, Jian Xie, Shen Wang, Zichao Wang, and Qingsong Wen. 2025. Uniedu: A unified language and vision assistant for education applications. *arXiv preprint arXiv:2503.20701*.
- Wildemakes de Almeida da Silva, Luis Carlos Costa Fonseca, Sofiane Labidi, and José Chrystian Lima Pacheco. 2024. Mitigation of hallucinations in language models in education: A new approach of comparative and cross-verification. In *2024 IEEE International Conference on Advanced Learning Technologies (ICALT)*, pages 207–209. IEEE.

| | | |
|------|--|------|
| 1025 | Shoutao Guo, Shaolei Zhang, Zhengrui Ma, Min Zhang, and Yang Feng. 2024b. Sillm: Large language models for simultaneous machine translation. <i>arXiv preprint arXiv:2402.13036</i> . | 1082 |
| 1026 | | 1083 |
| 1027 | | 1084 |
| 1028 | | 1085 |
| 1029 | Shuchen Guo, Ehsan Latif, Yifan Zhou, Xuan Huang, and Xiaoming Zhai. 2024c. Using generative ai and multi-agents to provide automatic feedback. <i>arXiv preprint arXiv:2411.07407</i> . | 1086 |
| 1030 | | 1087 |
| 1031 | | |
| 1032 | | 1088 |
| 1033 | Taicheng Guo, Bozhao Nan, Zhenwen Liang, Zhichun Guo, Nitesh Chawla, Olaf Wiest, Xiangliang Zhang, and 1 others. 2023. What can large language models do in chemistry? a comprehensive benchmark on eight tasks. <i>Advances in Neural Information Processing Systems</i> , 36:59662–59688. | 1089 |
| 1034 | | 1090 |
| 1035 | | 1091 |
| 1036 | | 1092 |
| 1037 | | 1093 |
| 1038 | | |
| 1039 | Jieun Han, Haneul Yoo, Junho Myung, Minsun Kim, Hyunseung Lim, Yoonsu Kim, Tak Yeon Lee, Hwajung Hong, Juho Kim, So-Yeon Ahn, and 1 others. 2023. Fabric: Automated scoring and feedback generation for essays. <i>arXiv preprint arXiv:2310.05191</i> . | 1094 |
| 1040 | | 1095 |
| 1041 | | 1096 |
| 1042 | | 1097 |
| 1043 | | 1098 |
| 1044 | Feng He, Tianqing Zhu, Dayong Ye, Bo Liu, Wanlei Zhou, and Philip S Yu. 2024a. The emerged security and privacy of llm agent: A survey with case studies. <i>arXiv preprint arXiv:2407.19354</i> . | 1099 |
| 1045 | | 1100 |
| 1046 | | |
| 1047 | | 1101 |
| 1048 | Zexue He, Yu Wang, An Yan, Yao Liu, Eric Chang, Amilcare Gentili, Julian McAuley, and Chun-Nan Hsu. 2023. MedEval: A multi-level, multi-task, and multi-domain medical benchmark for language model evaluation . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 8725–8744, Singapore. Association for Computational Linguistics. | 1102 |
| 1049 | | 1103 |
| 1050 | | 1104 |
| 1051 | | 1105 |
| 1052 | | 1106 |
| 1053 | | |
| 1054 | | 1107 |
| 1055 | | 1108 |
| 1056 | Zhitao He, Pengfei Cao, Chenhao Wang, Zhuoran Jin, Yubo Chen, Jiexin Xu, Huaijun Li, Xiaojian Jiang, Kang Liu, and Jun Zhao. 2024b. Agentscourt: Building judicial decision-making agents with court debate simulation and legal knowledge augmentation. <i>arXiv preprint arXiv:2403.02959</i> . | 1109 |
| 1057 | | 1110 |
| 1058 | | 1111 |
| 1059 | | |
| 1060 | | 1112 |
| 1061 | | 1113 |
| 1062 | Huu-Tuong Ho, Duc-Tin Ly, and Luong Vuong Nguyen. 2024. Mitigating hallucinations in large language models for educational application. In <i>2024 IEEE International Conference on Consumer Electronics-Asia (ICCE-Asia)</i> , pages 1–4. IEEE. | 1114 |
| 1063 | | 1115 |
| 1064 | | 1116 |
| 1065 | | |
| 1066 | | 1117 |
| 1067 | Shengxin Hong, Chang Cai, Sixuan Du, Haiyue Feng, Siyuan Liu, and Xiuyi Fan. 2024. " my grade is wrong!": A contestable ai framework for interactive feedback in evaluating student essays. <i>arXiv preprint arXiv:2409.07453</i> . | 1118 |
| 1068 | | 1119 |
| 1069 | | 1120 |
| 1070 | | 1121 |
| 1071 | | |
| 1072 | Jinchang Hou, Chang Ao, Haihong Wu, Xiangtao Kong, Zhigang Zheng, Daijia Tang, Chengming Li, Xiping Hu, Ruifeng Xu, Shiwen Ni, and 1 others. 2024. E-eval: a comprehensive chinese k-12 education evaluation benchmark for large language models. <i>arXiv preprint arXiv:2401.15927</i> . | 1122 |
| 1073 | | 1123 |
| 1074 | | 1124 |
| 1075 | | |
| 1076 | | 1125 |
| 1077 | | 1126 |
| 1078 | Bihao Hu, Jiayi Zhu, Yiyang Pei, and Xiaoqing Gu. 2025a. Exploring the potential of llm to enhance teaching plans through teaching simulation. <i>npj Science of Learning</i> , 10(1):7. | 1127 |
| 1079 | | 1128 |
| 1080 | | 1129 |
| 1081 | | 1130 |
| | | 1131 |
| | | 1132 |
| | | 1133 |
| | | |
| | Silan Hu and Xiaoning Wang. 2024. Foke: A personalized and explainable education framework integrating foundation models, knowledge graphs, and prompt engineering. In <i>China National Conference on Big Data and Social Computing</i> , pages 399–411. Springer. | 1134 |
| | | 1135 |
| | | 1136 |
| | | 1137 |
| | | 1138 |
| | | |
| | Yutao Hu, Tianbin Li, Quanfeng Lu, Wenqi Shao, Junjun He, Yu Qiao, and Ping Luo. 2024. Omnimedvqa: A new large-scale comprehensive evaluation benchmark for medical lvlm. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 22170–22183. | 1139 |
| | | 1140 |
| | | 1141 |
| | | 1142 |
| | | 1143 |
| | | |
| | Zhe Hu, Hou Pong Chan, Jing Li, and Yu Yin. 2025b. Debate-to-write: A persona-driven multi-agent framework for diverse argument generation . In <i>Proceedings of the 31st International Conference on Computational Linguistics</i> , pages 4689–4703, Abu Dhabi, UAE. Association for Computational Linguistics. | 1144 |
| | | 1145 |
| | | 1146 |
| | | 1147 |
| | | 1148 |
| | | 1149 |
| | | 1150 |
| | | |
| | Wenyue Hua, Xianjun Yang, Mingyu Jin, Zelong Li, Wei Cheng, Ruixiang Tang, and Yongfeng Zhang. 2024. Trustagent: Towards safe and trustworthy llm-based agents through agent constitution. In <i>Trustworthy Multi-modal Foundation Models and AI Agents (TiFA)</i> . | 1151 |
| | | 1152 |
| | | 1153 |
| | | 1154 |
| | | 1155 |
| | | 1156 |
| | | |
| | Dong Huang, Jie M Zhang, Michael Luck, Qingwen Bu, Yuhao Qing, and Heming Cui. 2023a. Agent-coder: Multi-agent-based code generation with iterative testing and optimisation. <i>arXiv preprint arXiv:2312.13010</i> . | 1157 |
| | | 1158 |
| | | 1159 |
| | | 1160 |
| | | 1161 |
| | | |
| | Haojing Huang, Jingheng Ye, Qingyu Zhou, Yinghui Li, Yangning Li, Feng Zhou, and Hai-Tao Zheng. 2023b. A frustratingly easy plug-and-play detection-and-reasoning module for chinese spelling check. <i>arXiv preprint arXiv:2310.09119</i> . | 1162 |
| | | 1163 |
| | | 1164 |
| | | 1165 |
| | | 1166 |
| | | |
| | Huazhen Huang, Xianguo Shi, Hongyang Lei, Fan Hu, and Yunpeng Cai. 2024a. Protchat: An ai multi-agent for automated protein analysis leveraging gpt-4 and protein language model. <i>Journal of Chemical Information and Modeling</i> , 65(1):62–70. | 1167 |
| | | 1168 |
| | | 1169 |
| | | 1170 |
| | | 1171 |
| | | |
| | Lan Huang. 2023. Ethics of artificial intelligence in education: Student privacy and data protection. <i>Science Insights Education Frontiers</i> , 16(2):2577–2587. | 1172 |
| | | 1173 |
| | | 1174 |
| | | |
| | Qian Huang, Jian Vora, Percy Liang, and Jure Leskovec. 2023c. Mlagentbench: Evaluating language agents on machine learning experimentation. <i>arXiv preprint arXiv:2310.03302</i> . | 1175 |
| | | 1176 |
| | | 1177 |
| | | 1178 |
| | | 1179 |
| | | 1180 |
| | | 1181 |
| | | |
| | Xu Huang, Weiwen Liu, Xiaolong Chen, Xingmei Wang, Hao Wang, Defu Lian, Yasheng Wang, Ruiming Tang, and Enhong Chen. 2024b. Understanding the planning of llm agents: A survey. <i>arXiv preprint arXiv:2402.02716</i> . | 1182 |
| | | 1183 |
| | | 1184 |
| | | 1185 |
| | | 1186 |
| | | 1187 |
| | | |
| | Jiahao Huo, Yibo Yan, Xu Zheng, Yuanhuiyi Lyu, Xin Zou, Zhihua Wei, and Xuming Hu. 2025. Mmunlearner: Reformulating multimodal machine unlearning in the era of multimodal large language models. <i>arXiv preprint arXiv:2502.11051</i> . | 1188 |
| | | 1189 |
| | | 1190 |
| | | 1191 |
| | | 1192 |
| | | 1193 |

| | | | |
|------|--|---|------|
| 1139 | Md Ashraful Islam, Mohammed Eunus Ali, and | Anupam Khan and Soumya K Ghosh. 2021. Student | 1194 |
| 1140 | Md Rizwan Parvez. 2024. Mapcoder: Multi-agent | performance analysis and prediction in classroom | 1195 |
| 1141 | code generation for competitive problem solving. | learning: A review of educational data mining studies. | 1196 |
| 1142 | <i>arXiv preprint arXiv:2405.11403</i> . | <i>Education and information technologies</i> , 26(1):205– | 1197 |
| | | 240. | 1198 |
| 1143 | Islam Asim Ismail. 2025. Protecting privacy in ai- | Wajahat Naseeb Khan. 2024. Ethical challenges of ai | 1199 |
| 1144 | enhanced education: A comprehensive examination | in education: Balancing innovation with data privacy. | 1200 |
| 1145 | of data privacy concerns and solutions in ai-based | <i>Journal of AI in Education: Innovations, Opportuni-</i> | 1201 |
| 1146 | learning. <i>Impacts of Generative AI on the Future of</i> | <i>ties, Challenges, and Future Directions</i> , 1(1):1–13. | 1202 |
| 1147 | <i>Research and Education</i> , pages 117–142. | | |
| 1148 | Kevin Maik Jablonka, Philippe Schwaller, Andres | Jing Yu Koh, Robert Lo, Lawrence Jang, Vikram | 1203 |
| 1149 | Ortega-Guerrero, and Berend Smit. 2024. Lever- | Duvvur, Ming Chong Lim, Po-Yu Huang, Graham | 1204 |
| 1150 | aging large language models for predictive chemistry. | Neubig, Shuyan Zhou, Ruslan Salakhutdinov, and | 1205 |
| 1151 | <i>Nature Machine Intelligence</i> , 6(2):161–169. | Daniel Fried. 2024. Visualwebarena: Evaluating mul- | 1206 |
| | | timodal agents on realistic visual web tasks. <i>arXiv</i> | 1207 |
| 1152 | Peter Jansen, Marc-Alexandre Côté, Tushar Khot, | <i>preprint arXiv:2401.13649</i> . | 1208 |
| 1153 | Erin Bransom, Bhavana Dalvi Mishra, Bod- | | |
| 1154 | hisattwa Prasad Majumder, Oyvind Taffjord, and Peter | Patrick Tser Jern Kon, Jiachen Liu, Qiuyi Ding, Yiming | 1209 |
| 1155 | Clark. 2024. Discoveryworld: A virtual environment | Qiu, Zhenning Yang, Yibo Huang, Jayanth Srinivasa, | 1210 |
| 1156 | for developing and evaluating automated scientific | Myungjin Lee, Mosharaf Chowdhury, and Ang Chen. | 1211 |
| 1157 | discovery agents. <i>Advances in Neural Information</i> | 2025. Curie: Toward rigorous and automated scien- | 1212 |
| 1158 | <i>Processing Systems</i> , 37:10088–10116. | tific experimentation with ai agents. <i>arXiv preprint</i> | 1213 |
| | | <i>arXiv:2502.16069</i> . | 1214 |
| 1159 | Hunkoog Jho. 2024. Leveraging generative ai in physics | Gerd Kortemeyer. 2023. Could an artificial-intelligence | 1215 |
| 1160 | education: Addressing hallucination issues in large | agent pass an introductory physics course? <i>Physical</i> | 1216 |
| 1161 | language models. | <i>Review Physics Education Research</i> , 19(1):010132. | 1217 |
| 1162 | Cong Jiang and Xiaolei Yang. 2024. Agents on the | Tomaž Kosar, Dragana Ostojić, Yu David Liu, and Mar- | 1218 |
| 1163 | bench: Large language model based multi agent | jan Mernik. 2024. Computer science education in | 1219 |
| 1164 | framework for trustworthy digital justice. <i>arXiv</i> | chatgpt era: Experiences from an experiment in a | 1220 |
| 1165 | <i>preprint arXiv:2412.18697</i> . | programming course for novice programmers. <i>Math-</i> | 1221 |
| 1166 | Zhoumingju Jiang and Mengjun Jiang. 2024. Be- | <i>ematics</i> , 12(5):629. | 1222 |
| 1167 | yond answers: Large language model-powered tu- | | |
| 1168 | toring system in physics education for deep learn- | Roman Koshkin, Katsuhito Sudoh, and Satoshi Naka- | 1223 |
| 1169 | ing and precise understanding. <i>arXiv preprint</i> | mura. 2024. Transllama: Llm-based simultaneous | 1224 |
| 1170 | <i>arXiv:2406.10934</i> . | translation system. <i>arXiv preprint arXiv:2402.04636</i> . | 1225 |
| 1171 | Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, | Lars Krupp, Steffen Steinert, Maximilian Kiefer- | 1226 |
| 1172 | Hanyi Fang, and Peter Szolovits. 2021. What disease | Emmanouilidis, Karina E Avila, Paul Lukowicz, | 1227 |
| 1173 | does this patient have? a large-scale open domain | Jochen Kuhn, Stefan Küchemann, and Jakob Karolus. | 1228 |
| 1174 | question answering dataset from medical exams. <i>Ap-</i> | 2024. Challenges and opportunities of moderating | 1229 |
| 1175 | <i>plied Sciences</i> , 11(14):6421. | usage of large language models in education. In <i>Pro-</i> | 1230 |
| 1176 | Hyoungwook Jin, Seonghee Lee, Hyungyu Shin, and | <i>ceedings of the International Conference on Mobile</i> | 1231 |
| 1177 | Juho Kim. 2024a. Teach ai how to code: Using | <i>and Ubiquitous Multimedia</i> , pages 249–254. | 1232 |
| 1178 | large language models as teachable agents for pro- | | |
| 1179 | gramming education. In <i>Proceedings of the 2024</i> | Max Ku, Thomas Chong, Jonathan Leung, Krish Shah, | 1233 |
| 1180 | <i>CHI Conference on Human Factors in Computing</i> | Alvin Yu, and Wenhui Chen. 2025. Theoremexplaina- | 1234 |
| 1181 | <i>Systems</i> , pages 1–28. | gent: Towards multimodal explanations for llm theo- | 1235 |
| 1182 | Hyoungwook Jin, Minju Yoo, Jeongeon Park, Yokyung | rem understanding . <i>Preprint</i> , arXiv:2502.19400. | 1236 |
| 1183 | Lee, Xu Wang, and Juho Kim. 2024b. Teachtune: | Kristjan-Julius Laak and Jaan Aru. 2024. Ai and per- | 1237 |
| 1184 | Reviewing pedagogical agents against diverse stu- | sonalized learning: bridging the gap with modern | 1238 |
| 1185 | dent profiles with simulated students. <i>arXiv preprint</i> | educational goals. <i>arXiv preprint arXiv:2404.02798</i> . | 1239 |
| 1186 | <i>arXiv:2410.04078</i> . | | |
| 1187 | Shi Jinxin, Zhao Jiabao, Wang Yilei, Wu Xingjiao, Li Ji- | Paraskevas Lagakis and Stavros Demetriadis. 2024. | 1240 |
| 1188 | awen, and He Liang. 2023. Cgmi: Configurable | Evaai: a multi-agent framework leveraging large lan- | 1241 |
| 1189 | general multi-agent interaction framework. <i>arXiv</i> | guage models for enhanced automated grading. In | 1242 |
| 1190 | <i>preprint arXiv:2308.12503</i> . | <i>International Conference on Intelligent Tutoring Sys-</i> | 1243 |
| | | <i>tems</i> , pages 378–385. Springer. | 1244 |
| 1191 | Mert Karabacak and Konstantinos Margetis. 2023. Em- | Jinqi Lai, Wensheng Gan, Jiayang Wu, Zhenlian Qi, and | 1245 |
| 1192 | bracing large language models for medical applica- | S Yu Philip. 2024. Large language models in law: A | 1246 |
| 1193 | tions: opportunities and challenges. <i>Cureus</i> , 15(5). | survey. <i>AI Open</i> . | 1247 |

| | | |
|------|---|------|
| 1248 | Jakub Lála, Odhran O'Donoghue, Aleksandar Shtedritski, Sam Cox, Samuel G Rodrigues, and Andrew D White. 2023. Paperqa: Retrieval-augmented generative agent for scientific research. <i>arXiv preprint arXiv:2312.07559</i> . | 1303 |
| 1249 | | 1304 |
| 1250 | | 1305 |
| 1251 | | 1306 |
| 1252 | | 1307 |
| 1253 | Soohwan Lee and Ki-Sang Song. 2024. Teachers' and students' perceptions of ai-generated concept explanations: Implications for integrating generative ai in computer science education. <i>Computers and Education: Artificial Intelligence</i> , 7:100283. | 1308 |
| 1254 | | 1309 |
| 1255 | | 1310 |
| 1256 | | 1311 |
| 1257 | | 1312 |
| 1258 | Bin Lei, Yi Zhang, Shan Zuo, Ali Payani, and Caiwen Ding. 2025. Macm: Utilizing a multi-agent system for condition mining in solving complex mathematical problems. <i>Advances in Neural Information Processing Systems</i> , 37:53418–53437. | 1313 |
| 1259 | | 1314 |
| 1260 | | 1315 |
| 1261 | | 1316 |
| 1262 | | 1317 |
| 1263 | Haitao Li, Junjie Chen, Jingli Yang, Qingyao Ai, Wei Jia, Youfeng Liu, Kai Lin, Yueyue Wu, Guozhi Yuan, Yiran Hu, and 1 others. 2024a. Legalagentbench: Evaluating llm agents in legal domain. <i>arXiv preprint arXiv:2412.17259</i> . | 1318 |
| 1264 | | 1319 |
| 1265 | | 1320 |
| 1266 | | 1321 |
| 1267 | | 1322 |
| 1268 | | 1323 |
| 1269 | | 1324 |
| 1270 | | 1325 |
| 1271 | | 1326 |
| 1272 | | 1327 |
| 1273 | | 1328 |
| 1274 | | 1329 |
| 1275 | | 1330 |
| 1276 | | 1331 |
| 1277 | | 1332 |
| 1278 | | 1333 |
| 1279 | | 1334 |
| 1280 | | 1335 |
| 1281 | | 1336 |
| 1282 | | 1337 |
| 1283 | | 1338 |
| 1284 | | 1339 |
| 1285 | | 1340 |
| 1286 | | 1341 |
| 1287 | | 1342 |
| 1288 | | 1343 |
| 1289 | | 1344 |
| 1290 | | 1345 |
| 1291 | | 1346 |
| 1292 | | 1347 |
| 1293 | | 1348 |
| 1294 | | 1349 |
| 1295 | | 1350 |
| 1296 | | 1351 |
| 1297 | | 1352 |
| 1298 | | 1353 |
| 1299 | | 1354 |
| 1300 | | 1355 |
| 1301 | | 1356 |
| 1302 | | 1357 |
| | | 1358 |
| | | 1359 |

| | | |
|------|--|------|
| 1360 | Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. 2024. The ai scientist: Towards fully automated open-ended scientific discovery. <i>arXiv preprint arXiv:2408.06292</i> . | 1416 |
| 1361 | | 1417 |
| 1362 | | |
| 1363 | | |
| 1364 | Rui Lv, Qi Liu, Weibo Gao, Haotian Zhang, Junyu Lu, and Linbo Zhu. 2025. Genal: Generative agent for adaptive learning. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 39, pages 577–585. | 1418 |
| 1365 | | 1419 |
| 1366 | | 1420 |
| 1367 | | 1421 |
| 1368 | | |
| 1369 | Andres M. Bran, Sam Cox, Oliver Schilter, Carlo Baldassari, Andrew D White, and Philippe Schwaller. 2024. Augmenting large language models with chemistry tools. <i>Nature Machine Intelligence</i> , 6(5):525–535. | 1422 |
| 1370 | | 1423 |
| 1371 | | 1424 |
| 1372 | | 1425 |
| 1373 | | 1426 |
| 1374 | | |
| 1375 | Pingchuan Ma, Tsun-Hsuan Wang, Minghao Guo, Zhiqing Sun, Joshua B Tenenbaum, Daniela Rus, Chuang Gan, and Wojciech Matusik. 2024a. Llm and simulation as bilevel optimizers: A new paradigm to advance physical scientific discovery. <i>arXiv preprint arXiv:2405.09783</i> . | 1427 |
| 1376 | | 1428 |
| 1377 | | 1429 |
| 1378 | | 1430 |
| 1379 | | |
| 1380 | Qianou Ma, Hua Shen, Kenneth Koedinger, and Sherry Tongshuang Wu. 2024b. How to teach programming in the ai era? using llms as a teachable agent for debugging. In <i>International Conference on Artificial Intelligence in Education</i> , pages 265–279. Springer. | 1431 |
| 1381 | | 1432 |
| 1382 | | 1433 |
| 1383 | | 1434 |
| 1384 | | 1435 |
| 1385 | | 1436 |
| 1386 | Yubo Ma, Zhibin Gou, Junheng Hao, Ruochen Xu, Shuohang Wang, Liangming Pan, Yujiu Yang, Yixin Cao, Aixin Sun, Hany Awadalla, and 1 others. 2024c. Sciagent: Tool-augmented language models for scientific reasoning. <i>arXiv preprint arXiv:2402.11451</i> . | 1437 |
| 1387 | | 1438 |
| 1388 | | 1439 |
| 1389 | | 1440 |
| 1390 | | 1441 |
| 1391 | Jakub Macina, Nico Daheim, Sankalan Pal Chowdhury, Tanmay Sinha, Manu Kapur, Iryna Gurevych, and Mrinmaya Sachan. 2023. Mathdial: A dialogue tutoring dataset with rich pedagogical properties grounded in math reasoning problems. <i>arXiv preprint arXiv:2305.14536</i> . | 1442 |
| 1392 | | 1443 |
| 1393 | | 1444 |
| 1394 | | 1445 |
| 1395 | | 1446 |
| 1396 | | 1447 |
| 1397 | Jakub Macina, Nico Daheim, Ido Hakimi, Manu Kapur, Iryna Gurevych, and Mrinmaya Sachan. 2025. Math-tutorbench: A benchmark for measuring open-ended pedagogical capabilities of llm tutors. <i>arXiv preprint arXiv:2502.18940</i> . | 1448 |
| 1398 | | 1449 |
| 1399 | | 1450 |
| 1400 | | 1451 |
| 1401 | | |
| 1402 | Bodhisattwa Prasad Majumder, Harshit Surana, Dhruv Agarwal, Bhavana Dalvi Mishra, Abhijeetsingh Meena, Aryan Prakhar, Tirth Vora, Tushar Khot, Ashish Sabharwal, and Peter Clark. 2024. Discoverybench: Towards data-driven discovery with large language models. <i>arXiv preprint arXiv:2407.01725</i> . | 1452 |
| 1403 | | 1453 |
| 1404 | | 1454 |
| 1405 | | 1455 |
| 1406 | | 1456 |
| 1407 | | 1457 |
| 1408 | Jordan K Matelsky, Felipe Parodi, Tony Liu, Richard D Lange, and Konrad P Kording. 2023. A large language model-assisted education tool to provide feedback on open-ended responses. <i>arXiv preprint arXiv:2308.02439</i> . | 1458 |
| 1409 | | 1459 |
| 1410 | | 1460 |
| 1411 | | 1461 |
| 1412 | | 1462 |
| 1413 | Siddharth Mehrotra, Carolina Centeio Jorge, Catholijn M Jonker, and Myrthe L Tielman. 2024. Integrity-based explanations for fostering appropriate trust in ai agents. <i>ACM Transactions on Interactive Intelligent Systems</i> , 14(1):1–36. | 1463 |
| 1414 | | 1464 |
| 1415 | | 1465 |
| | Silvia Milano, Joshua A McGrane, and Sabina Leonelli. 2023. Large language models challenge the future of higher education. <i>Nature Machine Intelligence</i> , 5(4):333–334. | 1466 |
| | | 1467 |
| | Paul Mithun, Enrique Noriega-Atala, Nirav Merchant, and Edwin Skidmore. 2025. Ai-verde: A gateway for egalitarian access to large language model-based resources for educational institutions. <i>arXiv preprint arXiv:2502.09651</i> . | 1468 |
| | | 1469 |
| | | 1470 |
| | Arindam Mitra, Hamed Khanpour, Corby Rosset, and Ahmed Awadallah. 2024. Orca-math: Unlocking the potential of slms in grade school math. <i>arXiv preprint arXiv:2402.14830</i> . | 1471 |
| | | 1472 |
| | Kaijie Mo and Renfen Hu. 2024. <i>ExpertEase: A multi-agent framework for grade-specific document simplification with large language models</i> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 9080–9099, Miami, Florida, USA. Association for Computational Linguistics. | 1473 |
| | | 1474 |
| | Hyeonseok Moon, Jaewook Lee, Sugyeong Eo, Chanjun Park, Jaehyung Seo, and Heui-Seok Lim. 2024. Generative interpretation: Toward human-like evaluation for educational question-answer pair generation. In <i>Findings of the Association for Computational Linguistics: EACL 2024</i> , pages 2185–2196. | 1475 |
| | | 1476 |
| | Michael Moret, Irene Pachon Angona, Leandro Cotos, Shen Yan, Kenneth Atz, Cyrill Brunner, Martin Baumgartner, Francesca Grisoni, and Gisbert Schneider. 2023. Leveraging molecular structure and bioactivity with chemical language models for de novo drug design. <i>Nature Communications</i> , 14(1):114. | 1477 |
| | | 1478 |
| | Christopher E Mower and Haitham Bou-Ammar. 2025. Al-khwarizmi: Discovering physical laws with foundation models. <i>arXiv preprint arXiv:2502.01702</i> . | 1479 |
| | | 1480 |
| | Inderjeet Nair, Jiaye Tan, Xiaotian Su, Anne Gere, Xu Wang, and Lu Wang. 2024. Closing the loop: Learning to generate writing feedback via language model simulated student revisions. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 16636–16657. | 1481 |
| | | 1482 |
| | Siddharth Narayanan, James D Braza, Ryan-Rhys Griffiths, Manu Ponnampati, Albert Bou, Jon Laurent, Ori Kabeli, Geemi Wellawatte, Sam Cox, Samuel G Rodrigues, and 1 others. 2024. Aviary: training language agents on challenging scientific tasks. <i>arXiv preprint arXiv:2412.21154</i> . | 1483 |
| | | 1484 |
| | Deepak Nathani, Lovish Madaan, Nicholas Roberts, Nikolay Bashlykov, Ajay Menon, Vincent Moens, Amar Budhiraja, Despoina Magka, Vladislav Vorotilov, Gaurav Chaurasia, and 1 others. 2025. Mlgym: A new framework and benchmark for advancing ai research agents. <i>arXiv preprint arXiv:2502.14499</i> . | 1485 |
| | | 1486 |
| | | 1487 |

| | | |
|------|---|------|
| 1471 | Jack Nelson. 2024. The other'Ilm': Large language models and the future of legal education. <i>European Journal of Legal Education</i> , 5(1):127–155. | 1528 |
| 1472 | | 1529 |
| 1473 | | 1530 |
| 1474 | Davy Tsz Kit Ng, Chee Wei Tan, and Jac Ka Lok Leung. 2024. Empowering student self-regulated learning and science education through chatgpt: A pioneering pilot study. <i>British Journal of Educational Technology</i> , 55(4):1328–1353. | 1531 |
| 1475 | | 1532 |
| 1476 | | 1533 |
| 1477 | | 1534 |
| 1478 | | |
| 1479 | Mengxu Pan, Alexandra Kitson, Hongyu Wan, and Mirjana Prpa. 2024. Ellma-t: an embodied llm-agent for supporting english language learning in social vr. <i>arXiv preprint arXiv:2410.02406</i> . | 1535 |
| 1480 | | 1536 |
| 1481 | | 1537 |
| 1482 | | 1538 |
| 1483 | Xinyu Pang, Ruixin Hong, Zhanke Zhou, Fangrui Lv, Xinwei Yang, Zhilong Liang, Bo Han, and Changshui Zhang. 2024. Physics reasoner: Knowledge-augmented reasoning for solving physics problems with large language models. <i>arXiv preprint arXiv:2412.13791</i> . | 1539 |
| 1484 | | |
| 1485 | | 1540 |
| 1486 | | 1541 |
| 1487 | | 1542 |
| 1488 | | 1543 |
| 1489 | Zachary A. Pardos, Matthew Tang, Ioannis Anastasopoulos, Shreya K. Sheel, and Ethan Zhang. 2023. Oatutor: An open-source adaptive tutoring system and curated content library for learning sciences research. In <i>Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems</i> , CHI '23, New York, NY, USA. Association for Computing Machinery. | 1544 |
| 1490 | | 1545 |
| 1491 | | 1546 |
| 1492 | | |
| 1493 | | 1547 |
| 1494 | | 1548 |
| 1495 | | 1549 |
| 1496 | | 1550 |
| 1497 | Minju Park, Sojung Kim, Seunghyun Lee, Soonwoo Kwon, and Kyuseok Kim. 2024. Empowering personalized learning through a conversation-based tutoring system with student modeling. In <i>Extended Abstracts of the CHI Conference on Human Factors in Computing Systems</i> , pages 1–10. | 1551 |
| 1498 | | |
| 1499 | | 1552 |
| 1500 | | 1553 |
| 1501 | | 1554 |
| 1502 | | 1555 |
| 1503 | Francesc Pedro, Miguel Subosa, Axel Rivas, and Paula Valverde. 2019. Artificial intelligence in education: Challenges and opportunities for sustainable development. | 1556 |
| 1504 | | |
| 1505 | | 1557 |
| 1506 | | 1558 |
| 1507 | Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, and 1 others. 2023. Toolllm: Facilitating large language models to master 16000+ real-world apis. <i>arXiv preprint arXiv:2307.16789</i> . | 1559 |
| 1508 | | 1560 |
| 1509 | | 1561 |
| 1510 | | |
| 1511 | | 1562 |
| 1512 | Nishat Raihan, Mohammed Latif Siddiq, Joanna CS Santos, and Marcos Zampieri. 2025. Large language models in computer science education: A systematic literature review. In <i>Proceedings of the 56th ACM Technical Symposium on Computer Science Education V. 1</i> , pages 938–944. | 1563 |
| 1513 | | 1564 |
| 1514 | | 1565 |
| 1515 | | 1566 |
| 1516 | | |
| 1517 | | 1567 |
| 1518 | Mayk Caldas Ramos, Christopher J Collison, and Andrew D White. 2025. A review of large language models and autonomous agents in chemistry. <i>Chemical Science</i> . | 1568 |
| 1519 | | 1569 |
| 1520 | | 1570 |
| 1521 | | 1571 |
| 1522 | Mahefa Abel Razafinirina, William Germain Dimbisoa, and Thomas Mahatody. 2024. Pedagogical alignment of large language models (llm) for personalized learning: a survey, trends and challenges. <i>Journal of Intelligent Learning Systems and Applications</i> , 16(4):448–480. | 1572 |
| 1523 | | 1573 |
| 1524 | | 1574 |
| 1525 | | |
| 1526 | | 1575 |
| 1527 | | 1576 |
| | | 1577 |
| | | 1578 |
| | | 1579 |
| | | |
| | | 1580 |
| | | 1581 |
| | | 1582 |
| | | 1583 |
| | | 1584 |
| | | |
| | | 1585 |
| | | 1586 |
| | | 1587 |
| | | 1588 |
| | | 1589 |
| | | 1590 |
| | | 1591 |
| | | 1592 |
| | | 1593 |
| | | 1594 |
| | | 1595 |
| | | 1596 |
| | | 1597 |
| | | 1598 |
| | | 1599 |
| | | 1600 |

| | | | |
|------|--|--|------|
| 1585 | Shuzheng Si, Wentao Ma, Haoyu Gao, Yuchuan Wu, | Ian Steenstra, Prasanth Murali, Rebecca B Perkins, | 1639 |
| 1586 | Ting-En Lin, Yinpei Dai, Hangyu Li, Rui Yan, Fei | Natalie Joseph, Michael K Paasche-Orlow, and Timothy | 1640 |
| 1587 | Huang, and Yongbin Li. 2023. Spokenwoz: A large- | Bickmore. 2024. Engaging and entertaining adoles- | 1641 |
| 1588 | scale speech-text benchmark for spoken task-oriented | cents in health education using llm-generated fantasy | 1642 |
| 1589 | dialogue agents. <i>Advances in Neural Information</i> | narrative games and virtual agents. In <i>Extended Ab-</i> | 1643 |
| 1590 | <i>Processing Systems</i> , 36:39088–39118. | <i>stracts of the CHI Conference on Human Factors in</i> | 1644 |
| | | <i>Computing Systems</i> , pages 1–8. | 1645 |
| 1591 | Nisha Simon and Christian Muise. 2022. Tattletale: | Jiamin Su, Yibo Yan, Fangteng Fu, Han Zhang, | 1646 |
| 1592 | storytelling with planning and large language models. | Jingheng Ye, Xiang Liu, Jiahao Huo, Huiyu Zhou, | 1647 |
| 1593 | In <i>ICAPS Workshop on Scheduling and Planning</i> | and Xuming Hu. 2025. Essayjudge: A multi-granular | 1648 |
| 1594 | <i>Applications</i> . | benchmark for assessing automated essay scoring | 1649 |
| 1595 | Michael D Skarlinski, Sam Cox, Jon M Laurent, | capabilities of multimodal large language models. | 1650 |
| 1596 | James D Braza, Michaela Hinks, Michael J Hammer- | <i>arXiv preprint arXiv:2502.11916</i> . | 1651 |
| 1597 | ling, Manvitha Ponnampati, Samuel G Rodriques, and | Theodore Summers, Shunyu Yao, Karthik Narasimhan, | 1652 |
| 1598 | Andrew D White. 2024. Language agents achieve | and Thomas Griffiths. Cognitive architectures for | 1653 |
| 1599 | superhuman synthesis of scientific knowledge. <i>arXiv</i> | language agents. <i>Transactions on Machine Learning</i> | 1654 |
| 1600 | <i>preprint arXiv:2409.13740</i> . | <i>Research</i> . | 1655 |
| 1601 | Milena Škobo and Vedran Petričević. 2023. Navigating | Jingyun Sun, Chengxiao Dai, Zhongze Luo, Yangbo | 1656 |
| 1602 | the challenges and opportunities of literary transla- | Chang, and Yang Li. 2024. Lawluo: A chinese | 1657 |
| 1603 | tion in the age of ai: Striking a balance between | law firm co-run by llm agents. <i>arXiv preprint</i> | 1658 |
| 1604 | human expertise and machine power. <i>Društvene i</i> | <i>arXiv:2407.16252</i> . | 1659 |
| 1605 | <i>humanističke studije</i> , 8(2 (23)):317–336. | | |
| 1606 | Samuel S Sohn, Danrui Li, Sen Zhang, Che-Jui Chang, | Melanie Swan, Takashi Kido, Eric Roland, and Renato | 1660 |
| 1607 | and Mubbasir Kapadia. 2024. From words to worlds: | P dos Santos. 2023. Math agents: Computational in- | 1661 |
| 1608 | Transforming one-line prompt into immersive multi- | frastructure, mathematical embedding, and genomics. | 1662 |
| 1609 | modal digital stories with communicative llm agent. | <i>arXiv preprint arXiv:2307.02502</i> . | 1663 |
| 1610 | <i>arXiv preprint arXiv:2406.10478</i> . | | |
| 1611 | Tao Song, Man Luo, Linjiang Chen, Yan Huang, Qing | Kehui Tan, Tianqi Pang, Chenyou Fan, and Song Yu. | 1664 |
| 1612 | Zhu, Daobin Liu, Baicheng Zhang, Gang Zou, Fei | 2023. Towards applying powerful large ai models in | 1665 |
| 1613 | Zhang, Weiwei Shang, and 1 others. 2024. A multi- | classroom teaching: Opportunities, challenges and | 1666 |
| 1614 | agent-driven robotic ai chemist enabling autonomous | prospects. <i>arXiv preprint arXiv:2305.03433</i> . | 1667 |
| 1615 | chemical research on demand. | | |
| 1616 | Xiangyu Song, Jianxin Li, Taotao Cai, Shuiqiao Yang, | Xiangru Tang, Yuliang Liu, Zefan Cai, Yanjun Shao, | 1668 |
| 1617 | Tingting Yang, and Chengfei Liu. 2022. A survey on | Junjie Lu, Yichi Zhang, Zexuan Deng, Helan Hu, | 1669 |
| 1618 | deep learning based knowledge tracing. <i>Knowledge-</i> | Kaikai An, Ruijun Huang, and 1 others. 2023. MI- | 1670 |
| 1619 | <i>Based Systems</i> , 258:110036. | bench: Evaluating large language models and agents | 1671 |
| | | for machine learning tasks on repository-level code. | 1672 |
| | | <i>arXiv preprint arXiv:2311.09835</i> . | 1673 |
| 1620 | Sinan Sonlu, Bennie Bendiksen, Funda Durupinar, | Minyang Tian, Luyu Gao, Shizhuo Zhang, Xinan Chen, | 1674 |
| 1621 | and Uğur Güdükbay. 2024. The effects of em- | Cunwei Fan, Xuefei Guo, Roland Haas, Pan Ji, Kit- | 1675 |
| 1622 | bodiment and personality expression on learning | tithat Krongchon, Yao Li, and 1 others. 2024a. Sci- | 1676 |
| 1623 | in llm-based educational agents. <i>arXiv preprint</i> | code: A research coding benchmark curated by sci- | 1677 |
| 1624 | <i>arXiv:2407.10993</i> . | entists. <i>Advances in Neural Information Processing</i> | 1678 |
| | | <i>Systems</i> , 37:30624–30650. | 1679 |
| 1625 | Henry W Sprueill, Carl Edwards, Khushbu Agarwal, | Shubo Tian, Qiao Jin, Lana Yeganova, Po-Ting Lai, | 1680 |
| 1626 | Mariefel V Olarte, Udishnu Sanyal, Conrad John- | Qingqing Zhu, Xiuying Chen, Yifan Yang, Qingyu | 1681 |
| 1627 | ston, Hongbin Liu, Heng Ji, and Sutanay Choud- | Chen, Won Kim, Donald C Comeau, and 1 others. | 1682 |
| 1628 | hury. 2024. Chemreasoner: Heuristic search over | 2024b. Opportunities and challenges for chatgpt and | 1683 |
| 1629 | a large language model’s knowledge space using | large language models in biomedicine and health. | 1684 |
| 1630 | quantum-chemical feedback. <i>arXiv preprint</i> | <i>Briefings in Bioinformatics</i> , 25(1):bbad493. | 1685 |
| 1631 | <i>arXiv:2402.10980</i> . | | |
| 1632 | Kamali N Sripathi, Rosa A Moscarella, Matthew Steele, | Meng-Lin Tsai, Chong Wei Ong, and Cheng-Liang | 1686 |
| 1633 | Rachel Yoho, Hyesun You, Luanna B Prevost, Mark | Chen. 2023. Exploring the use of large language | 1687 |
| 1634 | Urban-Lurain, John Merrill, and Kevin C Haudek. | models (llms) in chemical engineering education: | 1688 |
| 1635 | 2024. Machine learning mixed methods text analy- | Building core course problem models with chat-gpt. | 1689 |
| 1636 | sis: An illustration from automated scoring models | <i>Education for Chemical Engineers</i> , 44:71–95. | 1690 |
| 1637 | of student writing in biology education. <i>Journal of</i> | Ehsan Ullah, Anil Parwani, Mirza Mansoor Baig, and | 1691 |
| 1638 | <i>mixed methods research</i> , 18(1):48–70. | Rajendra Singh. 2024. Challenges and barriers of | 1692 |

| | | | |
|------|--|--|------|
| 1693 | using large language models (llm) such as chat- | Xiyu Wang, Yufei Wang, Satoshi Tsutsui, Weisi Lin, | 1751 |
| 1694 | gpt for diagnostic medicine with a focus on digi- | Bihan Wen, and Alex Kot. 2024e. Evolving story- | 1752 |
| 1695 | tal pathology—a recent scoping review. <i>Diagnostic</i> | telling: benchmarks and methods for new character | 1753 |
| 1696 | <i>pathology</i> , 19(1):43. | customization with diffusion models. In <i>Proceed-</i> | 1754 |
| | | <i>ings of the 32nd ACM International Conference on</i> | 1755 |
| 1697 | Karthik Valmeekam, Matthew Marquez, Alberto Olmo, | <i>Multimedia</i> , pages 3751–3760. | 1756 |
| 1698 | Sarath Sreedharan, and Subbarao Kambhampati. | | |
| 1699 | 2023. Planbench: An extensible benchmark for eval- | Yi Ru Wang, Jiafei Duan, Dieter Fox, and Siddhartha | 1757 |
| 1700 | uating large language models on planning and reason- | Srinivasa. 2023. Newton: Are large language mod- | 1758 |
| 1701 | ing about change. <i>Advances in Neural Information</i> | els capable of physical reasoning? <i>arXiv preprint</i> | 1759 |
| 1702 | <i>Processing Systems</i> , 36:38975–38987. | <i>arXiv:2310.07018</i> . | 1760 |
| | | | |
| 1703 | Guangya Wan, Yuqi Wu, Jie Chen, and Sheng Li. 2024. | Hao Wei, Jianing Qiu, Haibao Yu, and Wu Yuan. 2024. | 1761 |
| 1704 | Cot rerailer: Enhancing the reliability of large lan- | Medco: Medical education copilots based on a multi- | 1762 |
| 1705 | guage models in complex reasoning tasks through | agent framework. <i>arXiv preprint arXiv:2408.12496</i> . | 1763 |
| 1706 | error detection and correction. <i>arXiv preprint</i> | | |
| 1707 | <i>arXiv:2408.13940</i> . | Lilian Weng. 2023. Llm-powered autonomous agents. | 1764 |
| | | <i>lilianweng.github.io</i> . | 1765 |
| 1708 | Fei Wang, Qi Liu, Enhong Chen, Zhenya Huang, | Minghao Wu, Yulin Yuan, Gholamreza Haffari, and | 1766 |
| 1709 | Yu Yin, Shijin Wang, and Yu Su. 2022. Neuralcd: | Longyue Wang. 2024a. (perhaps) beyond human | 1767 |
| 1710 | a general framework for cognitive diagnosis. <i>IEEE</i> | translation: Harnessing multi-agent collaboration for | 1768 |
| 1711 | <i>Transactions on Knowledge and Data Engineering</i> , | translating ultra-long literary texts. <i>arXiv preprint</i> | 1769 |
| 1712 | 35(8):8312–8327. | <i>arXiv:2405.11804</i> . | 1770 |
| | | | |
| 1713 | Jian Wang, Yinpei Dai, Yichi Zhang, Ziqiao Ma, Wenjie | Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran | 1771 |
| 1714 | Li, and Joyce Chai. 2025a. Training turn-by-turn | Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun | 1772 |
| 1715 | verifiers for dialogue tutoring agents: The curious | Zhang, Shaokun Zhang, Jiale Liu, and 1 others. | 1773 |
| 1716 | case of llms as your coding tutors. <i>arXiv preprint</i> | 2023a. Autogen: Enabling next-gen llm applica- | 1774 |
| 1717 | <i>arXiv:2502.13311</i> . | tions via multi-agent conversation. <i>arXiv preprint</i> | 1775 |
| | | <i>arXiv:2308.08155</i> . | 1776 |
| 1718 | Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao | Xuansheng Wu, Haiyan Zhao, Yaochen Zhu, Yucheng | 1777 |
| 1719 | Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, | Shi, Fan Yang, Tianming Liu, Xiaoming Zhai, Wen- | 1778 |
| 1720 | Xu Chen, Yankai Lin, and 1 others. 2024a. A survey | lin Yao, Jundong Li, Mengnan Du, and 1 others. | 1779 |
| 1721 | on large language model based autonomous agents. | 2024b. Usable xai: 10 strategies towards exploit- | 1780 |
| 1722 | <i>Frontiers of Computer Science</i> , 18(6):186345. | ing explainability in the llm era. <i>arXiv preprint</i> | 1781 |
| | | <i>arXiv:2403.08946</i> . | 1782 |
| 1723 | Longyue Wang, Zefeng Du, Wenxiang Jiao, Chenyang | Yiran Wu. 2025. An empirical study on challenging | 1783 |
| 1724 | Lyu, Jianhui Pang, Leyang Cui, Kaiqiang Song, | math problem solving with llm-based conversational | 1784 |
| 1725 | Derek Wong, Shuming Shi, and Zhaopeng Tu. 2024b. | agents. | 1785 |
| 1726 | Benchmarking and improving long-text translation | | |
| 1727 | with large language models . In <i>Findings of the As-</i> | Yiran Wu, Feiran Jia, Shaokun Zhang, Hangyu Li, | 1786 |
| 1728 | <i>sociation for Computational Linguistics: ACL 2024</i> , | Erkang Zhu, Yue Wang, Yin Tat Lee, Richard Peng, | 1787 |
| 1729 | pages 7175–7187, Bangkok, Thailand. Association | Qingyun Wu, and Chi Wang. 2023b. Mathchat: Con- | 1788 |
| 1730 | for Computational Linguistics. | verse to tackle challenging math problems with llm | 1789 |
| | | agents. <i>arXiv preprint arXiv:2306.01337</i> . | 1790 |
| 1731 | Tianfu Wang, Yi Zhan, Jianxun Lian, Zhengyu Hu, | Chunqiu Steven Xia, Yinlin Deng, Soren Dunn, and | 1791 |
| 1732 | Nicholas Jing Yuan, Qi Zhang, Xing Xie, and Hui | Lingming Zhang. 2024. Agentless: Demystify- | 1792 |
| 1733 | Xiong. 2025b. Llm-powered multi-agent framework | ing llm-based software engineering agents. <i>arXiv</i> | 1793 |
| 1734 | for goal-oriented learning in intelligent tutoring sys- | <i>preprint arXiv:2407.01489</i> . | 1794 |
| 1735 | tem. <i>arXiv preprint arXiv:2501.15749</i> . | | |
| | | Nan Xie, Yuelin Bai, Hengyuan Gao, Ziqiang Xue, Feit- | 1795 |
| 1736 | Tiannan Wang, Jiamin Chen, Qingrui Jia, Shuai Wang, | eng Fang, Qixuan Zhao, Zhijian Li, Liang Zhu, Shi- | 1796 |
| 1737 | Ruoyu Fang, Huilin Wang, Zhaowei Gao, Chunzhao | wen Ni, and Min Yang. 2024a. Delilaw: A chinese | 1797 |
| 1738 | Xie, Chuou Xu, Jihong Dai, and 1 others. 2024c. | legal counselling system based on a large language | 1798 |
| 1739 | Weaver: Foundation models for creative writing. | model. In <i>Proceedings of the 33rd ACM Interna-</i> | 1799 |
| 1740 | <i>arXiv preprint arXiv:2401.17268</i> . | <i>tional Conference on Information and Knowledge</i> | 1800 |
| | | <i>Management</i> , pages 5299–5303. | 1801 |
| 1741 | Xidong Wang, Guiming Chen, Song Dingjie, Zhang | Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan | 1802 |
| 1742 | Zhiyi, Zhihong Chen, Qingying Xiao, Junying Chen, | Li, Siheng Zhao, Ruisheng Cao, Jing Hua Toh, Zhou- | 1803 |
| 1743 | Feng Jiang, Jianquan Li, Xiang Wan, Benyou Wang, | jun Cheng, Dongchan Shin, Fangyu Lei, and 1 others. | 1804 |
| 1744 | and Haizhou Li. 2024d. CMB: A comprehensive | 2024b. Osloworld: Benchmarking multimodal agents | 1805 |
| 1745 | medical benchmark in Chinese . In <i>Proceedings of</i> | | |
| 1746 | <i>the 2024 Conference of the North American Chap-</i> | | |
| 1747 | <i>ter of the Association for Computational Linguistics:</i> | | |
| 1748 | <i>Human Language Technologies (Volume 1: Long</i> | | |
| 1749 | <i>Papers)</i> , pages 6184–6205, Mexico City, Mexico. As- | | |
| 1750 | sociation for Computational Linguistics. | | |

| | | | |
|------|--|---|------|
| 1806 | for open-ended tasks in real computer environments. | John Yang, Carlos Jimenez, Alexander Wettig, Kilian | 1863 |
| 1807 | <i>Advances in Neural Information Processing Systems</i> , | Lieret, Shunyu Yao, Karthik Narasimhan, and Ofir | 1864 |
| 1808 | 37:52040–52094. | Press. 2025. Swe-agent: Agent-computer interfaces | 1865 |
| 1809 | Wei Xiong, Chengshuai Shi, Jiaming Shen, Aviv Rosen- | enable automated software engineering. <i>Advances in</i> | 1866 |
| 1810 | berg, Zhen Qin, Daniele Calandriello, Misha Khal- | <i>Neural Information Processing Systems</i> , 37:50528– | 1867 |
| 1811 | man, Rishabh Joshi, Bilal Piot, Mohammad Saleh, | 50652. | 1868 |
| 1812 | and 1 others. 2024. Building math agents with multi- | John Yang, Carlos E Jimenez, Alex L Zhang, Kil- | 1869 |
| 1813 | turn iterative preference learning. <i>arXiv preprint</i> | ian Lieret, Joyce Yang, Xindi Wu, Ori Press, | 1870 |
| 1814 | <i>arXiv:2409.02392</i> . | Niklas Muennighoff, Gabriel Synnaeve, Karthik R | 1871 |
| 1815 | Songlin Xu, Hao-Ning Wen, Hongyi Pan, Dallas | Narasimhan, and 1 others. 2024a. Swe-bench multi- | 1872 |
| 1816 | Dominguez, Dongyin Hu, and Xinyu Zhang. 2025. | modal: Do ai systems generalize to visual software | 1873 |
| 1817 | Classroom simulacra: Building contextual stu- | domains? <i>arXiv preprint arXiv:2410.03859</i> . | 1874 |
| 1818 | dent generative agents in online education for | Kaiqi Yang, Yucheng Chu, Taylor Darwin, Ahreum Han, | 1875 |
| 1819 | learning behavioral simulation. <i>arXiv preprint</i> | Hang Li, Hongzhi Wen, Yasemin Copur-Gencturk, | 1876 |
| 1820 | <i>arXiv:2502.02780</i> . | Jiliang Tang, and Hui Liu. 2024b. Content knowl- | 1877 |
| 1821 | Songlin Xu, Xinyu Zhang, and Lianhui Qin. 2024a. Ed- | edge identification with multi-agent large language | 1878 |
| 1822 | uagent: Generative student agents in learning. <i>arXiv</i> | models (llms). In <i>International Conference on Ar-</i> | 1879 |
| 1823 | <i>preprint arXiv:2404.07963</i> . | <i>tificial Intelligence in Education</i> , pages 284–292. | 1880 |
| 1824 | Tianlong Xu, Yi-Fan Zhang, Zhendong Chu, Shen | Springer. | 1881 |
| 1825 | Wang, and Qingsong Wen. 2024b. Ai-driven virtual | Ziqi Yang, Xuhai Xu, Bingsheng Yao, Ethan | 1882 |
| 1826 | teacher for enhanced educational efficiency: Leverag- | Rogers, Shao Zhang, Stephen Intille, Nawar Shara, | 1883 |
| 1827 | ing large pretrain models for autonomous error analy- | Guodong Gordon Gao, and Dakuo Wang. 2024c. | 1884 |
| 1828 | sis and correction. <i>arXiv preprint arXiv:2409.09403</i> . | Talk2care: An llm-based voice assistant for com- | 1885 |
| 1829 | Lixiang Yan, Lele Sha, Linxuan Zhao, Yuheng Li, | munication between healthcare providers and older | 1886 |
| 1830 | Roberto Martinez-Maldonado, Guanliang Chen, | adults. <i>Proceedings of the ACM on Interactive, Mo-</i> | 1887 |
| 1831 | Xinyu Li, Yueqiao Jin, and Dragan Gašević. 2024a. | <i>bile, Wearable and Ubiquitous Technologies</i> , 8(2):1– | 1888 |
| 1832 | Practical and ethical challenges of large language | 35. | 1889 |
| 1833 | models in education: A systematic scoping review. | Zonglin Yang, Wanhao Liu, Ben Gao, Tong Xie, | 1890 |
| 1834 | <i>British Journal of Educational Technology</i> , 55(1):90– | Yuqiang Li, Wanli Ouyang, Soujanya Poria, Erik | 1891 |
| 1835 | 112. | Cambria, and Dongzhan Zhou. 2024d. Moose- | 1892 |
| 1836 | Yibo Yan and Joey Lee. 2024. Georeasoner: Reason- | chem: Large language models for rediscovering un- | 1893 |
| 1837 | ing on geospatially grounded context for natural lan- | seen chemistry scientific hypotheses. <i>arXiv preprint</i> | 1894 |
| 1838 | guage understanding. In <i>Proceedings of the 33rd</i> | <i>arXiv:2410.07076</i> . | 1895 |
| 1839 | <i>ACM International Conference on Information and</i> | Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, | 1896 |
| 1840 | <i>Knowledge Management</i> , pages 4163–4167. | Tom Griffiths, Yuan Cao, and Karthik Narasimhan. | 1897 |
| 1841 | Yibo Yan, Jiamin Su, Jianxiang He, Fangteng Fu, | 2023. Tree of thoughts: Deliberate problem solving | 1898 |
| 1842 | Xu Zheng, Yuanhuiyi Lyu, Kun Wang, Shen Wang, | with large language models. <i>Advances in neural</i> | 1899 |
| 1843 | Qingsong Wen, and Xuming Hu. 2024b. A survey | <i>information processing systems</i> , 36:11809–11822. | 1900 |
| 1844 | of mathematical reasoning in the era of multimodal | Jingheng Ye, Yong Jiang, Xiaobin Wang, Yinghui Li, | 1901 |
| 1845 | large language model: Benchmark, method & chal- | Yangning Li, Hai-Tao Zheng, Pengjun Xie, and Fei | 1902 |
| 1846 | lenges. <i>arXiv preprint arXiv:2412.11936</i> . | Huang. 2024a. Productagent: Benchmarking conver- | 1903 |
| 1847 | Yibo Yan, Shen Wang, Jiahao Huo, Xuming Hu, and | sational product search agent with asking clarification | 1904 |
| 1848 | Qingsong Wen. 2025a. Mathagent: Leveraging | questions. <i>arXiv preprint arXiv:2407.00942</i> . | 1905 |
| 1849 | a mixture-of-math-agent framework for real-world | Jingheng Ye, Yinghui Li, Yangning Li, and Hai-Tao | 1906 |
| 1850 | multimodal mathematical error detection. <i>arXiv</i> . | Zheng. 2023a. Mixedit: Revisiting data augmen- | 1907 |
| 1851 | Yibo Yan, Shen Wang, Jiahao Huo, Hang Li, Boyan | tation and beyond for grammatical error correction. | 1908 |
| 1852 | Li, Jiamin Su, Xiong Gao, Yi-Fan Zhang, Tianlong | <i>arXiv preprint arXiv:2310.11671</i> . | 1909 |
| 1853 | Xu, Zhendong Chu, and 1 others. 2024c. Errorradar: | Jingheng Ye, Yinghui Li, Shirong Ma, Rui Xie, Wei | 1910 |
| 1854 | Benchmarking complex mathematical reasoning of | Wu, and Hai-Tao Zheng. 2022. Focus is what you | 1911 |
| 1855 | multimodal large language models via error detection. | need for chinese grammatical error correction. <i>arXiv</i> | 1912 |
| 1856 | <i>arXiv preprint arXiv:2410.04509</i> . | <i>preprint arXiv:2210.12692</i> . | 1913 |
| 1857 | Yibo Yan, Shen Wang, Jiahao Huo, Jingheng Ye, Zhen- | Jingheng Ye, Yinghui Li, and Haitao Zheng. 2023b. | 1914 |
| 1858 | dong Chu, Xuming Hu, Philip S Yu, Carla Gomes, | System report for ccl23-eval task 7: Thu kelab (sz)- | 1915 |
| 1859 | Bart Selman, and Qingsong Wen. 2025b. Posi- | exploring data augmentation and denoising for chi- | 1916 |
| 1860 | tion: Multimodal large language models can signifi- | nese grammatical error correction. In <i>Proceedings</i> | 1917 |
| 1861 | cantly advance scientific reasoning. <i>arXiv preprint</i> | <i>of the 22nd Chinese National Conference on Compu-</i> | 1918 |
| 1862 | <i>arXiv:2502.02871</i> . | <i>tational Linguistics (Volume 3: Evaluations)</i> , pages | 1919 |
| | | 262–270. | 1920 |

| | | | |
|------|--|--|------|
| 1921 | Jingheng Ye, Yinghui Li, Qingyu Zhou, Yangning Li, Shirong Ma, Hai-Tao Zheng, and Ying Shen. 2023c. Cleme: debiasing multi-reference evaluation for grammatical error correction. <i>arXiv preprint arXiv:2305.10819</i> . | Dong Zhang, Zhaowei Li, Pengyu Wang, Xin Zhang, Yaqian Zhou, and Xipeng Qiu. 2024a. Speechagents: Human-communication simulation with multi-modal multi-agent systems. <i>arXiv preprint arXiv:2401.03945</i> . | 1977 |
| 1922 | | | 1978 |
| 1923 | | | 1979 |
| 1924 | | | 1980 |
| 1925 | | | 1981 |
| 1926 | Jingheng Ye, Shang Qin, Yinghui Li, Xuxin Cheng, Libo Qin, Hai-Tao Zheng, Peng Xing, Zishan Xu, Guo Cheng, and Zhao Wei. 2024b. Excgec: A benchmark of edit-wise explainable chinese grammatical error correction. <i>arXiv preprint arXiv:2407.00924</i> . | Kechi Zhang, Jia Li, Ge Li, Xianjie Shi, and Zhi Jin. 2024b. Codeagent: Enhancing code generation with tool-integrated agent systems for real-world repo-level coding challenges. <i>arXiv preprint arXiv:2401.07339</i> . | 1982 |
| 1927 | | | 1983 |
| 1928 | | | 1984 |
| 1929 | | | 1985 |
| 1930 | | | 1986 |
| 1931 | Jingheng Ye, Shang Qin, Yinghui Li, Hai-Tao Zheng, Shen Wang, and Qingsong Wen. 2025a. Corrections meet explanations: A unified framework for explainable grammatical error correction. <i>arXiv preprint arXiv:2502.15261</i> . | Mike Zhang, Amalie Pernille Dilling, Léon Gondelman, Niels Erik Ruan Lyngdorf, Euan D Lindsay, and Johannes Bjerva. 2025a. Seft: Harnessing large language model agents to improve educational feedback systems. <i>arXiv preprint arXiv:2502.12927</i> . | 1987 |
| 1932 | | | 1988 |
| 1933 | | | 1989 |
| 1934 | | | 1990 |
| 1935 | | | 1991 |
| 1936 | Jingheng Ye, Shen Wang, Deqing Zou, Yibo Yan, Kun Wang, Hai-Tao Zheng, Zenglin Xu, Irwin King, Philip S Yu, and Qingsong Wen. 2025b. Position: Llms can be good tutors in foreign language education. <i>arXiv preprint arXiv:2502.05467</i> . | Rongsheng Zhang, Jiji Tang, Chuanqi Zang, Mingtao Pei, Wei Liang, Zeng Zhao, and Zhou Zhao. 2025b. Let storytelling tell vivid stories: A multi-modal-agent-based unified storytelling framework. <i>Neuro-computing</i> , 622:129316. | 1992 |
| 1937 | | | 1993 |
| 1938 | | | 1994 |
| 1939 | | | 1995 |
| 1940 | | | 1996 |
| 1941 | Jingheng Ye, Zishan Xu, Yinghui Li, Xuxin Cheng, Linlin Song, Qingyu Zhou, Hai-Tao Zheng, Ying Shen, and Xin Su. 2024c. Cleme2. 0: Towards more interpretable evaluation by disentangling edits for grammatical error correction. <i>arXiv preprint arXiv:2407.00934</i> . | Xueqiao Zhang, Chao Zhang, Jianwen Sun, Jun Xiao, Yi Yang, and Yawei Luo. 2025c. Eduplanner: Llm-based multi-agent systems for customized and intelligent instructional design. <i>IEEE Transactions on Learning Technologies</i> . | 1997 |
| 1942 | | | 1998 |
| 1943 | | | 1999 |
| 1944 | | | 2000 |
| 1945 | | | 2001 |
| 1946 | | | |
| 1947 | Botao Yu, Frazier N Baker, Ziru Chen, Garrett Herb, Boyu Gou, Daniel Adu-Ampratwum, Xia Ning, and Huan Sun. 2024. Tooling or not tooling? the impact of tools on language agents for chemistry problem solving. <i>arXiv preprint arXiv:2411.07228</i> . | Yi-Fan Zhang, Hang Li, Dingjie Song, Lichao Sun, Tianlong Xu, and Qingsong Wen. 2025d. From correctness to comprehension: Ai agents for personalized error diagnosis in education. <i>arXiv preprint arXiv:2502.13789</i> . | 2002 |
| 1948 | | | 2003 |
| 1949 | | | 2004 |
| 1950 | | | 2005 |
| 1951 | | | 2006 |
| 1952 | Weikang Yuan, Junjie Cao, Zhuoren Jiang, Yangyang Kang, Jun Lin, Kaisong Song, Pengwei Yan, Changlong Sun, Xiaozhong Liu, and 1 others. 2024. Can large language models grasp legal theories? enhance legal reasoning with insights from multi-agent collaboration. <i>arXiv preprint arXiv:2410.02507</i> . | Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, and 1 others. 2023. Siren's song in the ai ocean: A survey on hallucination in large language models. <i>arXiv preprint arXiv:2309.01219</i> , 2(5). | 2007 |
| 1953 | | | 2008 |
| 1954 | | | 2009 |
| 1955 | | | 2010 |
| 1956 | | | 2011 |
| 1957 | | | |
| 1958 | Murong Yue, Wenhan Lyu, Wijdan Mifdal, Jennifer Suh, Yixuan Zhang, and Ziyu Yao. 2024. Mathvc: An llm-simulated multi-character virtual classroom for mathematics education. <i>arXiv preprint arXiv:2404.06711</i> . | Zaibin Zhang, Yongting Zhang, Lijun Li, Hongzhi Gao, Lijun Wang, Huchuan Lu, Feng Zhao, Yu Qiao, and Jing Shao. 2024c. Psysafe: A comprehensive framework for psychological-based attack, defense, and evaluation of multi-agent system safety. <i>arXiv preprint arXiv:2401.11880</i> . | 2012 |
| 1959 | | | 2013 |
| 1960 | | | 2014 |
| 1961 | | | 2015 |
| 1962 | | | 2016 |
| 1963 | Osmar R Zaiane. 2002. Building a recommender agent for e-learning systems. In <i>International Conference on Computers in Education, 2002. Proceedings.</i> , pages 55–59. IEEE. | Zeyu Zhang, Xiaohe Bo, Chen Ma, Rui Li, Xu Chen, Quanyu Dai, Jieming Zhu, Zhenhua Dong, and Jirong Wen. 2024d. A survey on the memory mechanism of large language model based agents. <i>arXiv preprint arXiv:2404.13501</i> . | 2017 |
| 1964 | | | 2018 |
| 1965 | | | 2019 |
| 1966 | | | 2020 |
| 1967 | Siyu Zha, Yuehan Qiao, Qingyu Hu, Zhongsheng Li, Jiangtao Gong, and Yingqing Xu. 2024. Designing child-centric ai learning environments: Insights from llm-enhanced creative project-based learning. <i>arXiv preprint arXiv:2403.16159</i> . | Haiteng Zhao, Chang Ma, Fangzhi Xu, Lingpeng Kong, and Zhi-Hong Deng. 2025. Biomaze: Benchmarking and enhancing large language models for biological pathway reasoning. <i>arXiv preprint arXiv:2502.16660</i> . | 2021 |
| 1968 | | | 2022 |
| 1969 | | | 2023 |
| 1970 | | | 2024 |
| 1971 | | | 2025 |
| 1972 | Xuesong Zhai, Xiaoyan Chu, Ching Sing Chai, Morris Siu Yung Jong, Andreja Istenic, Michael Spector, Jia-Bao Liu, Jing Yuan, and Yan Li. 2021. A review of artificial intelligence (ai) in education from 2010 to 2020. <i>Complexity</i> , 2021(1):8812542. | Jiawei Zheng, Hanghai Hong, Feiyan Liu, Xiaoli Wang, Jingsong Su, Yonggui Liang, and Shikai Wu. 2024. Fine-tuning large language models for domain-specific machine translation. <i>arXiv preprint arXiv:2402.15061</i> . | 2026 |
| 1973 | | | 2027 |
| 1974 | | | 2028 |
| 1975 | | | 2029 |
| 1976 | | | 2030 |
| | | | 2031 |
| | | | 2032 |

- Longwei Zheng, Fei Jiang, Xiaoqing Gu, Yuanyuan Li, Gong Wang, and Haomin Zhang. 2025. Teaching via llm-enhanced simulations: Authenticity and barriers to suspension of disbelief. *The Internet and Higher Education*, 65:100990.
- Wanjun Zhong, Lianghong Guo, Qiqi Gao, and Yanlin Wang. 2023. Memorybank: Enhancing large language models with long-term memory. *arXiv preprint arXiv:2305.10250*.
- Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue Ou, Yonatan Bisk, Daniel Fried, and 1 others. 2023. Webarena: A realistic web environment for building autonomous agents. *arXiv preprint arXiv:2307.13854*.
- Xi Zhu, Yu Wang, Hang Gao, Wujiang Xu, Chen Wang, Zhiwei Liu, Kun Wang, Mingyu Jin, Linsey Pang, Qingsong Wen, and 1 others. 2024. Recommender systems meet large language model agents: A survey. *Available at SSRN 5062105*.
- Yuchen Zhuang, Yue Yu, Kuan Wang, Haotian Sun, and Chao Zhang. 2023. Toolqa: A dataset for llm question answering with external tools. *Advances in Neural Information Processing Systems*, 36:50117–50143.
- Deqing Zou, Jingheng Ye, Yulu Liu, Yu Wu, Zishan Xu, Yinghui Li, Hai-Tao Zheng, Bingxu An, Zhao Wei, and Yong Xu. 2025. Revisiting classification taxonomy for grammatical errors. *arXiv preprint arXiv:2502.11890*.

A Domain-Specific Educational Agents

Recent research on LLM agents in education has also shown growing interest in domain-specific applications. We explore their use in *science learning*, *language learning*, and *professional development*, focusing on their algorithmic frameworks, agentic designs, and relevant datasets and benchmarks. In Figure 3, we present a systematic taxonomy of domain-specific educational agents.

A.1 Agent for Science Learning

An agent for science learning is an intelligent system powered by LLMs, designed to assist students in acquiring and applying scientific knowledge through personalized, interactive experiences (Yan et al., 2025b; Raihan et al., 2025; Ng et al., 2024; Brown and Cruz Castro, 2025). The significance of these agents in education lies in their ability to offer tailored feedback, enhance conceptual understanding, and promote active engagement with complex scientific ideas. In the following sections, we explore the impact of LLM agents in four key scientific disciplines: *mathematics* (§A.1.1), *physics* (§A.1.2), *chemistry* (§A.1.3), and *biology* (§A.1.4), as well as their broader contributions to *general scientific discovery* (§A.1.5).

A.1.1 Mathematics

In mathematics, LLM agents provide substantial support by helping students navigate complex problems and reinforcing their understanding of abstract concepts (Yan et al., 2024b; Xiong et al., 2024; Swan et al., 2023; Yan et al., 2024c; Wu, 2025; Mitra et al., 2024). For instance, Gou et al. (2023) introduce TORA (Tool-integrated Reasoning Agents), a framework that integrates natural language reasoning and program-based tool use to handle mathematical reasoning. MathAgent (Yan et al., 2025a) similarly proposes Mixture-of-Math-Agent framework to address multimodal error detection in real-world K-12 scenarios, and flexibly transform the visual information of different types of questions into forms that are more easily understood by LLMs (e.g., converting plane geometry images into formalized expression). Additionally, MathChat (Wu et al., 2023b) serves as a conversational mathematical problem-solving agent, which consists of a chat-based LLM agent and a tool-based user agent. Furthermore, Xiong et al. (2024) propose to use reinforcement learning from human feedback (RLHF) to further improve tool-integrated

agents for mathematical problem-solving, and formulate this method as a Markov decision process, distinguishing it from the typical contextual bandit approach used in RLHF. Besides, MACM (Lei et al., 2025) discuss the limitations of LLMs in handling complex mathematical logical deduction, thus introducing a multi-agent system, which comprises three interactive agents: Thinker, Judge, and Executor.

A.1.2 Physics

In the field of physics, LLM agents help students make sense of challenging concepts and offer interactive tools to simulate physical phenomena (Pang et al., 2024; Mower and Bou-Ammar, 2025; Barman et al., 2025; Feng et al.; Jiang and Jiang, 2024; Yan and Lee, 2024). Wang et al. (2023) introduce NEWTON, the first pipeline and benchmark to explore the physical reasoning abilities of LLMs. Furthermore, Kortemeyer (2023) describe a case study exploring if an LLM agent can pass an introductory calculus-based physics course. In addition, Physics Reasoner (Pang et al., 2024), a novel knowledge-augmented framework for physics problem-solving, leverages a comprehensive formula set and detailed checklists to ensure accuracy and completeness. It can serve as an agent consisting of three stages - problem analysis, formula retrieval, and guided reasoning. Besides, Ma et al. (2024a) describe the Scientific Generative Agent (SGA), a bilevel optimization framework designed for physical scientific discovery, and highlight the use of LLMs for generating and revising scientific hypotheses and implementing an exploit-and-explore strategy.

A.1.3 Chemistry

Chemistry education also benefits greatly from LLM agents, which can explain molecular structures, chemical reactions, and experimental processes in an engaging and interactive way (Ramos et al., 2025; M. Bran et al., 2024; Yu et al., 2024; Guo et al., 2023; Tsai et al., 2023). For example, ChemCrow (M. Bran et al., 2024) is the first LLM chemistry agent capable of autonomous planning and execution of chemical syntheses, including an insect repellent and three organocatalysts. Yu et al. (2024) further present ChemAgent, an enhanced chemistry agent improved over ChemCrow, with a focus on two essential cognitive abilities of chemistry problem-solving: reasoning and grounding. Besides, Curie (Kon et al., 2025) is an agent framework aimed at incorporating rigor into the experi-

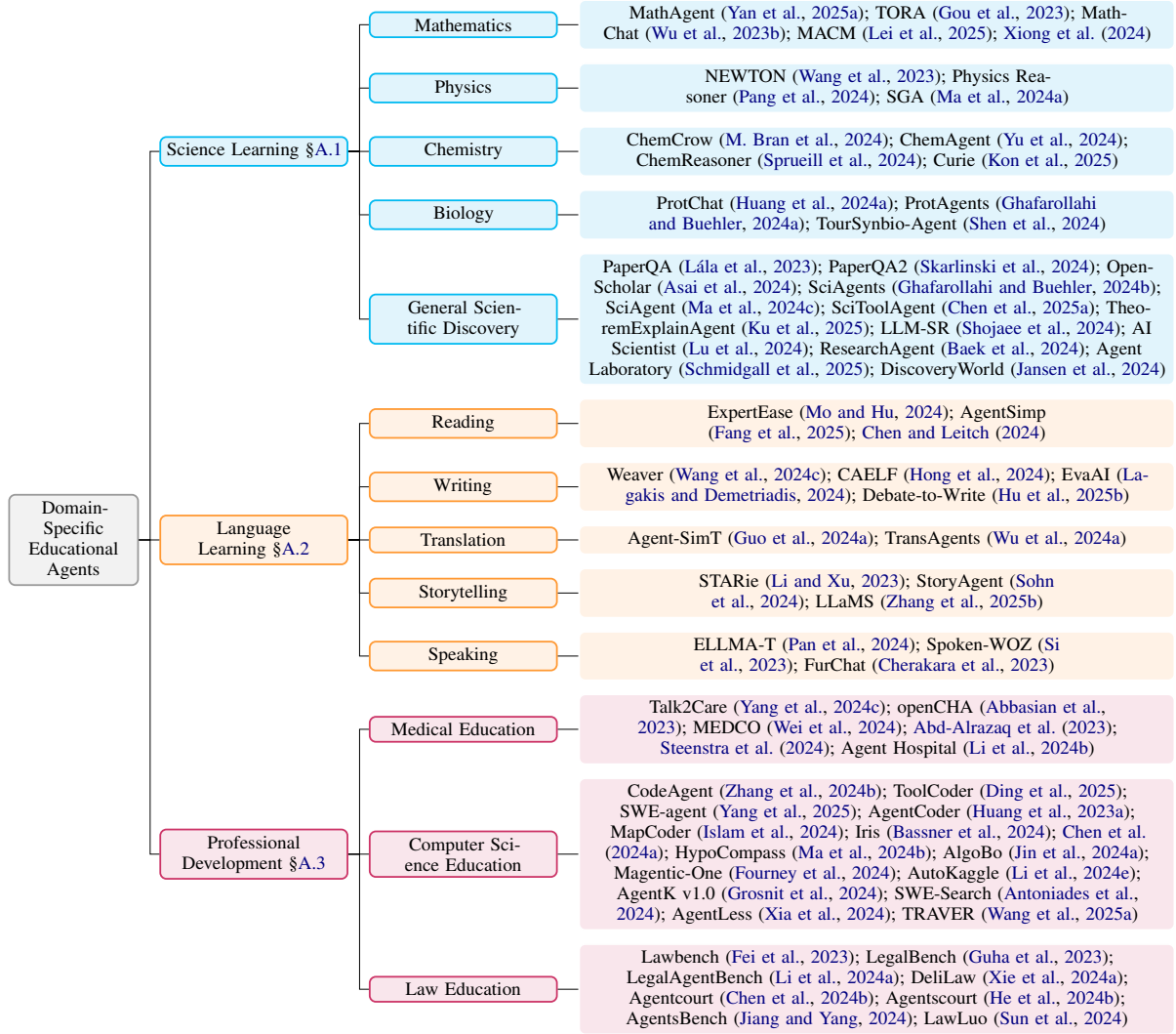


Figure 3: Taxonomy of domain-specific educational agents.

mentation process via three core elements: an intra-agent rigor module to boost reliability, an inter-agent rigor module to ensure systematic control, and an experiment knowledge module to improve interoperability. Recent studies have explored the capabilities of LLMs in complex chemical discovery (Yang et al., 2024d; Ruan et al., 2024; Sprueill et al., 2024; Moret et al., 2023; Jablonka et al., 2024), and their potential can be advanced by leveraging the interactivity of agent-based tool use and the flexibility of planning strategies (Song et al., 2024; Ramos et al., 2025).

A.1.4 Biology

In biology, LLM agents enhance learning by offering detailed explanations of biological processes and providing interactive experiences to explore living systems (Yan et al., 2025b; Bhattacharya et al., 2024; Sripathi et al., 2024; Zhao et al., 2025; Gao et al., 2024). For example, ProtChat (Huang et al., 2024a) is a multi-agent tool leveraging GPT-4 and Protein Language Models for seamless protein analysis automation, thus evolutionizing the complexities of protein sequence interpretation. ProtAgents (Ghafarollahi and Buehler, 2024a) is introduced as a multi-agent modeling framework that combines state-of-the-art LLMs with diverse tools to tackle protein design and analysis. It consists of a team of agents: User, Planner, Assistant, Critic, and Group Chat Manager. Besides, Shen et al. (2024) present TourSynbio-Agent, an innovative agent framework that leverages TourSynbio-7B’s protein understanding ability to perform various protein engineering tasks, such as mutation analysis, inverse folding, and visualization.

A.1.5 General Scientific Discovery

LLM agents support general scientific discovery by assisting students in data interpretation, hypothesis testing, and creative problem-solving (Yan et al., 2025b; Chen et al., 2024d; Ghafarollahi and Buehler, 2024b; Chen et al., 2025a; Schmidgall et al., 2025). These LLM agents, such as PaperQA (Lála et al., 2023), PaperQA2 (Skarliniski et al., 2024), OpenScholar (Asai et al., 2024), SciAgents (Ghafarollahi and Buehler, 2024b), TheoremExplainAgent (Ku et al., 2025), and LLM-SR (Shojaee et al., 2024), can analyze complex scientific datasets, helping students uncover patterns and trends that may not be immediately apparent. In addition, Narayanan et al. (2024) present Aviary, an extensible gymnasium for language agents for

three challenging scientific tasks: manipulating DNA constructs for molecular cloning, answering research questions by accessing scientific literature, and engineering protein stability. Furthermore, both SciAgent (Ma et al., 2024c) and SciToolAgent (Chen et al., 2025a) extend to a tool-augmented scientific reasoning setting with the help of domain-specific tools. Besides, Agent Laboratory (Schmidgall et al., 2025) emerges as an agent framework that automates the research process of three phases (Literature Review, Experimentation, and Report Writing) via various LLM agents (PhD, Postdoc, ML Engineer, *etc.*).

A.2 Agent for Language Learning

The integration of LLM agents into language learning is revolutionizing how core competencies—reading, writing, listening, and speaking—are taught and practiced (Ye et al., 2025b). These skills form the foundation of effective communication and language acquisition, and recent advancements in LLM-based agents have significantly enhanced how learners interact with and acquire these skills. Below, we introduce recent advancements in each subdomain, highlighting the role of LLM agents in enhancing pedagogical outcomes in language acquisition for students and second language (L2) speakers through engaging and adaptive approaches (Huang et al., 2023b; Ye et al., 2023b,a).

A.2.1 Reading

Reading comprehension is a vital component of language learning, and LLM agents are playing an increasingly important role in enhancing students’ reading abilities. For instance, ExpertEase (Mo and Hu, 2024) employs a multi-agent framework to adapt documents for grade-specific audiences, simulating expert-teacher-student collaboration to enhance comprehension. AgentSimp (Fang et al., 2025) tackles document-level simplification by leveraging multiple agents with distinct roles to ensure coherence and accessibility. Additionally, various LLMs (Adetayo et al., 2024) have been used as academic reading companions, demonstrating improved engagement and understanding of complex qualitative texts in educational settings (Chen and Leitch, 2024).

A.2.2 Writing

The development of writing skills has benefited significantly from NLP tasks like explainable gram-

mathematical error correction (EXGEC) (Ye et al., 2024b, 2025a; Zou et al., 2025) and automatic essay scoring (AES) systems (Su et al., 2025). Weaver (Wang et al., 2024c), a family of LLMs fine-tuned for writing tasks (Ye et al., 2023c, 2024c), outperforms generalist LLMs like GPT-4 in generating human-like narratives. Moreover, Weaver natively supports retrieval-augmented generation (RAG) and function calling, serving as a qualified foundational model for LLM agents. For interactive feedback on student essays, CAELF (Hong et al., 2024) introduces a multi-agent framework that enables interactive essay feedback. By combining Teaching-Assistant agents’ evaluations with teacher-agent arbitration, students can contest grades and engage with the feedback, addressing the “black box” limitations of traditional automated scoring. Inspired by the process of human debate, Debate-to-Write (Hu et al., 2025b) construct a persona-based multi-agent framework that can enable agents to collaboratively debate, discuss ideas, and form a comprehensive plan for argument writing.

A.2.3 Translation

LLM agents demonstrate remarkable advancements in both simultaneous (Koshkin et al., 2024; Guo et al., 2024b) and literary translation (Cheng et al., 2024; Škobo and Petričević, 2023). Translation tasks benefit from LLM agents through their ability to integrate specialized tools and orchestrate multi-agent collaboration. Agent-SiMT (Guo et al., 2024a) combines the decision-making capabilities of a Simultaneous Machine Translation (SiMT) policy agent with the generative power of a translation agent, achieving state-of-the-art performance in simultaneous translation by dynamically balancing reading and generation actions. For literary translation, TransAgents (Wu et al., 2024a) employs a multi-agent framework to replicate the complex workflows of human translation teams, addressing cultural nuances and stylistic challenges through collaborative reasoning. This approach not only improves translation quality but also extends LLM applications to linguistically and culturally rich domains. These contributions underscore the importance of tool use and agent collaboration in advancing translation education (Zheng et al., 2024).

A.2.4 Storytelling

Storytelling applications leverage LLM agents to create immersive and interactive learning experi-

ences (Simon and Muise, 2022). STARie (Li and Xu, 2023), a peer-like embodied conversational agent, integrates multimodal tools such as speech synthesis and facial animation to scaffold children’s storytelling, fostering narrative creativity and oral communication skills (Beredo and Ong, 2021; Cas-sell, 2022). StoryAgent (Sohn et al., 2024) combines top-down story drafting with bottom-up asset generation to transform simple prompts into coherent, multi-modal digital narratives. By automating complex storytelling workflows (Liem et al., 2023), it democratizes content creation and enhances engagement in language learning. LLaMS (Zhang et al., 2025b), a multi-modal agent framework, is designed to generate multi-modal human-level stories characterized by expressiveness and consistency, incorporating the Story-Adapter module for long image sequence illustration. These systems demonstrate the potential of LLM agents to support both cognitive and creative aspects of language education for children.

A.2.5 Speaking

LLM agents are revolutionizing spoken language education by integrating reasoning and multi-agent collaboration to build adaptive dialogue systems (Liu et al., 2024d; Balan et al., 2024). ELLMA-T (Pan et al., 2024) employs contextual reasoning and role-playing in social VR environments to provide personalized feedback and language assessments, enabling learners to practice speaking in realistic scenarios (Lim et al., 2024; Li et al., 2025b). SpokenWOZ (Si et al., 2023) introduces a large-scale benchmark for task-oriented spoken dialogue, highlighting the importance of reasoning and multi-turn interaction in addressing real-world conversational challenges. FurChat (Cherakara et al., 2023), an embodied conversational agent, combines verbal and non-verbal communication cues to simulate natural interactions, making it a valuable tool for improving speaking skills through immersive and realistic practice. By employing multimodal signals such as speech and gestures, SpeechAgents (Zhang et al., 2024a) enhances the authenticity of dialogue simulations, capturing consistent content, natural rhythm, and rich emotional expression. Through Multi-Agent Tuning (Liang et al., 2024), it optimizes LLM capabilities for large-scale simulations involving up to 25 agents, enabling applications like drama creation and audio novel generation.

A.3 Agent for Professional Development

Agents for professional development harness the capabilities of LLMs to offer scalable, adaptive, and context-aware learning experiences tailored to domain-specific needs. This section summarizes how recent studies develop agents to revolutionize professional training in fields including *medical* (§A.3.1), *computer science* (§A.3.2), and *law education* (§A.3.3).

A.3.1 Medical Education

The deployment of LLM agents in healthcare has created new opportunities for personalized, interactive, and scalable systems, with several health agents introduced (Shusterman et al., 2025) such as Talk2Care (Yang et al., 2024c) and openCHA (Abbasian et al., 2023). Additionally, Abd-Alrazaq et al. (2023) highlight the educational potentials of LLMs in crafting personalized curricula, adaptive learning plans, and dynamic assessment tools for medical education, while concurrently addressing challenges including algorithmic bias, misinformation, and privacy issues. MEDCO (Wei et al., 2024), a multi-agent system, has the capacity to replicate real-world medical training environments through agent collaboration with virtual patients, expert physicians, and radiologists, enhancing interdisciplinary learning and peer interaction. Furthermore, Abbasian et al. (2023) introduce openCHA as a personalized LLM-powered framework that integrates external resources and orchestrates multi-step problem-solving for complex healthcare queries (Ye et al., 2024a), emphasizing tool use and action planning. Beyond traditional education, Steenstra et al. (2024) explore LLMs in creating fantasy narrative games for adolescent health education, demonstrating the agents’ ability to generate engaging, doctor-validated content that enhances knowledge retention through gamification. Li et al. (2024b) present Agent Hospital, a simulation environment where LLM-driven agents evolve through autonomous interactions, demonstrating significant improvements in medical reasoning and performance on benchmarks like MedQA (Jin et al., 2021) after treating thousands of simulated patients. Collectively, these investigations highlight the versatility of LLM agents within medical education, demonstrating their abilities in reasoning, collaboration, tool integration, and adaptive learning to effectively address a broad spectrum of educational and clinical challenges (Karabacak and Margetis, 2023; Tian et al., 2024b; Ullah et al., 2024).

A.3.2 Computer Science Education

An agent for computer science (CS) education greatly enhances learning by providing personalized guidance on coding, debugging, and understanding CS principles (Ma et al., 2024b; Lee and Song, 2024; Kosar et al., 2024; Liu et al., 2024b). For example, CodeAgent (Zhang et al., 2024b) serves as an LLM agent framework for repo-level code generation, incorporating external tools such as WebSearch and DocSearch. Recent studies have demonstrated the potential of agent-based code generation systems such as ToolCoder (Ding et al., 2025), SWE-agent (Yang et al., 2025), AgentCoder (Huang et al., 2023a), and MapCoder (Islam et al., 2024), which can significantly enhance students’ coding efficiency (Jin et al., 2024a; Wang et al., 2025a; Frankford et al., 2024). Furthermore, Bassner et al. (2024) introduce Iris, an LLM-driven virtual tutor designed to offer personalized, context-aware assistance to CS students within the interactive learning platform Artemis. Besides, Chen et al. (2024a) propose Learning-by-Teaching (LBT) as an effective pedagogical strategy for CS education, and leverage the advantages of LLM agents (e.g., contextual conversation & learning from demonstrations).

A.3.3 Law Education

LLM agents leverage pre-trained legal knowledge, interactive capabilities, and reasoning skills to support law education through judicial interpretation, moot court simulation, and case analysis (Chen et al., 2024b; Nelson, 2024; Lai et al., 2024; Yuan et al., 2024). However, evaluations from LawBench (Fei et al., 2023) and LegalBench (Guha et al., 2023) reveal that LLMs struggle with legal knowledge application and judicial aid. LegalAgentBench (Li et al., 2024a) further highlights their limitations in multi-hop reasoning and defense statement writing, showing that LLM agents require significant improvements to effectively assist in complex legal tasks. Despite these challenges, LLM agents are emerging as valuable tools for moot court simulations, a crucial component of legal reasoning and advocacy training. DeliLaw (Xie et al., 2024a) enhances law education by integrating legal and case retrieval modules, enabling students to practice legal research, case analysis, statutory interpretation, and mock consultations. LawLuo (Sun et al., 2024) applies a multi-agent framework with retrieval-augmented generation to simulate multi-turn legal consultations, improving

personalization and ambiguity handling. Similarly, AgentCourt (Chen et al., 2024b) and AgentsCourt (He et al., 2024b) simulate courtroom interactions and judicial decision-making, providing a realistic training ground for law students. AgentsBench (Jiang and Yang, 2024) extends this by offering multi-agent legal reasoning and case analysis, further advancing AI-driven legal education.

B Datasets & Benchmarks

In Table 2, we provide a comprehensive summary of publicly available datasets and benchmarks designed to evaluate LLM agents for education across various domains. It categorizes resources based on their primary goal, target users, subject domain, education level, language, modality and dataset size. We hope this collection can support and advance research on LLM agents for education.

Several datasets are designed to evaluate the pedagogical agents, such as ASSIST09 (Feng et al., 2009) and Junyi (Chang et al., 2015), which support knowledge tracing (KT) in K-12 math education, while others like EduAgent (Xu et al., 2024a) facilitate adaptive learning (AL) by dynamically adjusting content based on student profiles. In addition, error correction and detection (ECD) datasets, such as Virtual Teacher (Xu et al., 2024b) and MathCCS (Zhang et al., 2025d), assess LLM agents’ ability to identify and rectify student mistakes in math learning. Other datasets cater to writing, reading, and language learning, including FABRIC (Han et al., 2023), EssayJudge (Su et al., 2025), and EX-CGEC (Ye et al., 2024b), which focus on feedback and generation (FCG) for student essays. Multi-Sim (Ryan et al., 2023) and Wang et al. (2024b) provide multi-lingual translation and storytelling benchmarks, expanding LLM capabilities beyond English-language education.

Several datasets support domain-specific educational agents across science, law, medicine, and computer science. Beyond their primary goal of evaluating pedagogical ability, these datasets assess LLM agents in domain-specific applications. They provide insights into how LLM agents can be adapted for specialized instruction, evaluating their ability to deliver subject-specific knowledge, facilitate problem-solving, and enhance interactive learning experiences across diverse educational fields. ScienceAgentBench (Chen et al., 2024d) and TheoremExplainBench (Ku et al., 2025) assess scientific reasoning and theorem explanation,

while ML-Bench (Tang et al., 2023) and MLA-gentBench (Huang et al., 2023c) focus on machine learning education. In law education, datasets like LawBench (Fei et al., 2023), LegalBench (Guha et al., 2023), and AgentCourt (Chen et al., 2024b) evaluate legal knowledge application, case analysis, and court simulations. Medical education datasets, including MedBench (Cai et al., 2024) and OmniMedVQA (Hu et al., 2024), test clinical reasoning and medical knowledge retrieval. For computer science, SWE-Bench (Yang et al., 2025) and Programming Feedback (Estévez-Ayres et al., 2024) assess code generation, debugging, and software engineering instruction. These benchmarks help refine LLM agents for specialized tutoring, enhancing AI-driven learning in professional fields.

Table 2: Summary of existing datasets and benchmarks of LLM agents for education.

| Dataset&Benchmark | Goal | User | Domain | Level | Language | Modality | Amount | Source |
|----------------------|-----------|---------|-----------------------------------|----------|---------------|--------------|-----------|------------------------------|
| ASSIST09 | KT | Student | Math | K12 | EN | text | 227k | (Feng et al., 2009) |
| Junyi | KT | Student | Math | K12 | ZH | text | 2.5M | (Chang et al., 2015) |
| EduAgent | AL & CS | Student | - | Graduate | EN | text & image | 1,015 | (Xu et al., 2024a) |
| MathDial | AL & KT | Student | Math | K12 | EN | text | 45 | (Macina et al., 2023) |
| MultiArith | AL | Student | Math | K12 | EN | text | 180 | (Xu et al., 2024a) |
| CoMTA | KT | Student | Math | K12 | EN | text | 153 | (Scarlato et al., 2025) |
| MaCKT | KT | Student | Math | K12 | EN | text | 452 | (Yang et al., 2024b) |
| Virtual Teacher | ECD | Student | Math | K12 | ZH | text & image | 420 | (Xu et al., 2024b) |
| MathCCS | ECD | Student | Math | K12 | ZH | text & image | 420 | (Zhang et al., 2025d) |
| MathTutorBench | ECD & FCG | Student | Math | K12 | EN | text | 7 tasks | (Macina et al., 2025) |
| MultiSim | - | Student | Reading | - | Multi-lingual | text | 1.7M | (Ryan et al., 2023) |
| FABRIC | FCG | Student | Writing | - | EN | text | 1,782 | (Han et al., 2023) |
| EssayJudge | FCG | Student | Writing | - | EN | text & image | 1,054 | (Su et al., 2025) |
| PROF | FCG | Teacher | Writing | - | EN | text | 363 | (Nair et al., 2024) |
| EXCGEC | FCG | Student | Writing | - | ZH | text | 8,216 | (Ye et al., 2024b) |
| Wang et al. (2024b) | - | All | Translation | - | Multi-lingual | text | 70K | (Wang et al., 2024b) |
| NewEpisode | - | Student | Storytelling | - | EN | text & image | 24.5K | (Wang et al., 2024e) |
| SD-Eval | ECD | Student | Speaking | - | EN | speech | 7,303 | (Ao et al., 2024) |
| Programming Feedback | FCG | Teacher | Computer Science | - | Code | text | 52 | (Estévez-Ayres et al., 2024) |
| Review Critique | FCG | All | Computer Science | - | EN | text | 440 | (Du et al., 2024) |
| AAAR-1.0 | FCG | All | Computer Science | - | EN | text & image | 1,000 | (Lou et al., 2024) |
| ScienceAgentBench | - | All | General Science | - | EN | text | 102 | (Chen et al., 2024d) |
| TheoremExplainBench | - | All | General Science | - | EN | text | 240 | (Ku et al., 2025) |
| MLGym-Bench | - | All | General Science | - | EN | text | 13 tasks | (Nathani et al., 2025) |
| ML-Bench | - | All | Machine Learning | - | EN | text & image | 9,641 | (Tang et al., 2023) |
| MLAgentBench | - | All | Machine Learning | - | EN | text | 13 tasks | (Huang et al., 2023c) |
| SciCode | - | All | General Science | - | EN | text | 338 | (Tian et al., 2024a) |
| BLADE | - | All | General Science | - | EN | text | 12 tasks | (Gu et al., 2024) |
| DiscoveryBench | - | All | General Science | - | EN | text | 1,167 | (Majumder et al., 2024) |
| SUPER | - | All | General Science | - | EN | text | 796 | (Bogin et al., 2024) |
| E-EVAL | - | All | Math & Language & General Science | K12 | ZH | text | 4,351 | (Hou et al., 2024) |
| MedBench | - | All | Medical | - | ZH | text | 40,041 | (Cai et al., 2024) |
| CMB | - | All | Medical | - | ZH | text | 280,839 | (Wang et al., 2024d) |
| OmniMedVQA | - | All | Medical | - | ZH | text & image | 127,995 | (Hu et al., 2024) |
| MedEval | - | All | Medical | - | EN | text & image | 22,779 | (He et al., 2023) |
| OSWorld | - | All | Computer Science | - | EN | text & image | 369 | (Xie et al., 2024b) |
| Spider2-V | - | All | Computer Science | - | EN | text & image | 812 | (Cao et al., 2024) |
| VisualWebArena | - | All | Computer Science | - | EN | text & image | 910 | (Koh et al., 2024) |
| WebArena | - | All | Computer Science | - | EN | text & image | 812 | (Zhou et al., 2023) |
| SWE-Bench | - | All | Computer Science | - | EN | text | 2,294 | (Yang et al., 2025) |
| SWE-Bench M | - | All | Computer Science | - | EN | text & image | 617 | (Yang et al., 2024a) |
| Magentic-One | - | All | Computer Science | - | EN | text & image | 617 | (Fourney et al., 2024) |
| Lawbench | - | All | Law | - | ZH | text | 20 tasks | (Fei et al., 2023) |
| LegalBench | - | All | Law | - | EN | text | 162 tasks | (Guha et al., 2023) |
| LegalAgentBench | - | All | Law | - | ZH | text | 300 tasks | (Li et al., 2024a) |
| Agentcourt | - | All | Law | - | ZH | text | 550 | (Chen et al., 2024b) |
| SimuCourt | - | All | Law | - | ZH | text | 420 | (He et al., 2024b) |