

# Position: LLMs Can be Good Tutors in English Education

Anonymous ACL submission

## Abstract

While recent efforts have begun integrating large language models (LLMs) into English education, they often rely on traditional approaches to learning tasks without fully embracing educational methodologies, thus lacking adaptability to language learning. To address this gap, we argue that **LLMs have the potential to serve as effective tutors in English Education**. Specifically, LLMs can play three critical roles: (1) as *data enhancers*, improving the creation of learning materials or serving as student simulations; (2) as *task predictors*, serving as learner assessment or optimizing learning pathway; and (3) as *agents*, enabling personalized and inclusive education. We encourage interdisciplinary research to explore these roles, fostering innovation while addressing challenges and risks, ultimately advancing English Education through the thoughtful integration of LLMs.

## 1 Introduction

English Education has long been a cornerstone of global education and a critical component of K-12 curricula, equipping students with the linguistic and cultural competencies necessary for an interconnected world (Alhusaiyan, 2025; Katinskaia, 2025). However, traditional English teaching methods often fall short in addressing the diverse needs of learners (Hou, 2020). Challenges such as limited personalization, scalability constraints, and the lack of real-time feedback are particularly pronounced in large classroom settings (Ehrenberg et al., 2001). Addressing these shortcomings requires innovative approaches that not only enhance the quality of instruction but also adapt to the unique learning trajectories of students (Eaton, 2010).

Recently, the advent of LLMs has opened new possibilities for transforming English Education (Caines et al., 2023). LLMs exhibit remarkable natural language understanding and generation capabilities, making them promising candidates for

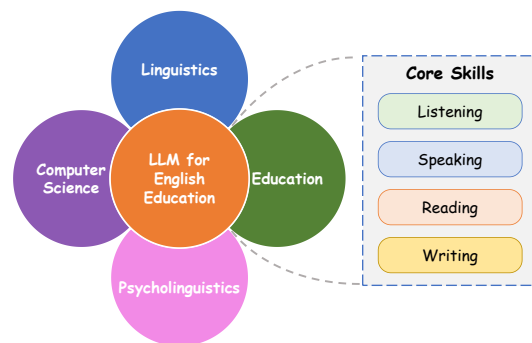


Figure 1: Involved disciplines of LLM for English Edu.

roles traditionally filled by human tutors. Leveraging LLMs as AI tutors can overcome inherent limitations of conventional teaching methods, offering scalable, interactive, and personalized learning experiences (Chen et al., 2024; Schmucker et al., 2024). Therefore, this position paper argues that **LLMs can be effective tutors in English education, complementing human expertise and addressing key limitations of traditional methods**.

As shown in Figure 1, English Education intersects with multiple *disciplines*, each of which underscores the potential of LLMs to revolutionize this domain. From the perspective of (1) *computer science*, advancements in machine learning and NLP have enabled LLMs to process and generate human-like language at an unprecedented scale; (2) *linguistics* (Radford et al., 2009) contributes a deeper understanding of grammar, phonetics, and semantics, allowing LLMs to generate accurate and understandable language outputs; (3) *education* provides the foundation for designing effective pedagogical strategies, ensuring that LLMs can deliver personalized, engaging, and developmentally appropriate learning experiences; and finally, (4) *psycholinguistics* (Steinberg and Sciarini, 2013) bridges the gap between language acquisition and cognitive processes, enabling LLMs to optimize learner interactions by adapting to individual needs and fostering meaningful engagement. Together,

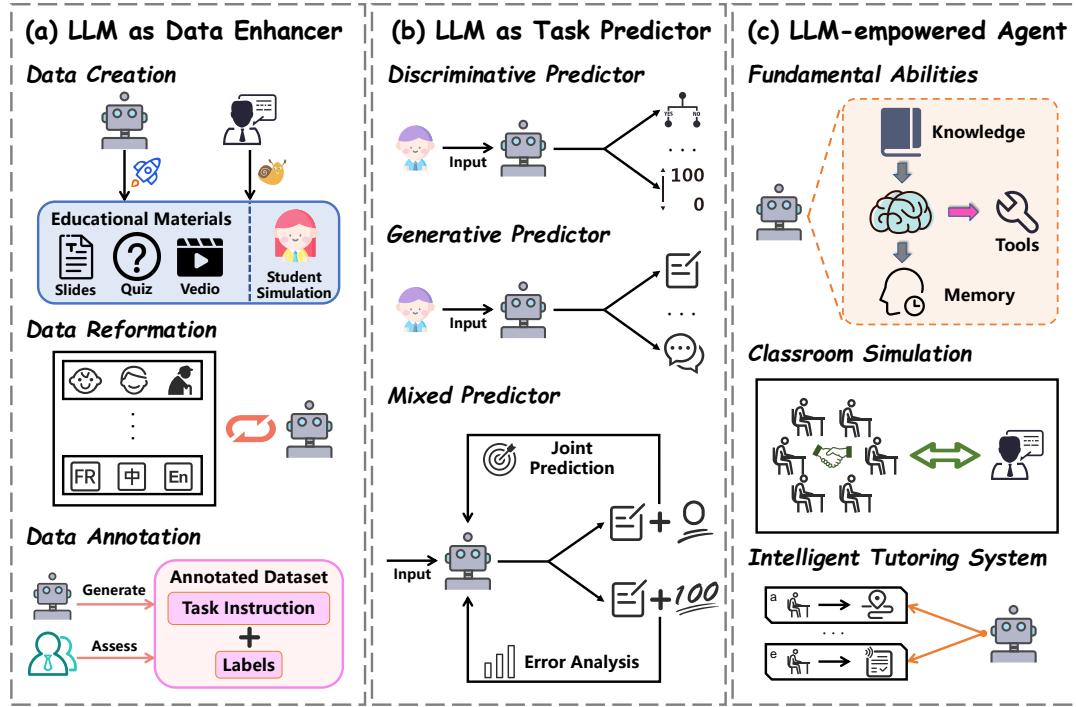


Figure 2: Overview of three roles of LLMs in English education. An overview of related literature is provided in Appendix A.

these disciplines position LLMs as uniquely capable of addressing the multifaceted challenges of English education.

Moreover, English Education encompasses four *core skills*: listening, speaking, reading, and writing, each of which can be significantly enhanced by LLMs. For listening, LLMs can generate diverse audio materials (Ghosal et al., 2023) and facilitate interactive voice-based exercises, helping learners improve their ability to discern pronunciation, intonation, and contextual meaning. In speaking, LLMs can simulate realistic conversations (Siyan et al., 2024), provide pronunciation feedback, and scaffold learners’ oral communication skills through iterative practice. For reading, LLMs can curate leveled texts, generate comprehension questions (Samuel et al., 2024), and engage learners in discussions that deepen their understanding of written content. Finally, in writing, LLMs can offer real-time grammar, syntax, and style feedback while assisting with idea generation and iterative revisions (Stahl et al., 2024a). By addressing these core skills holistically, LLMs have the potential to deliver a comprehensive and adaptive learning experience.

Despite these opportunities, the deployment of LLMs in English Education must be approached carefully, ensuring that their integration comple-

ments rather than replaces human tutors (Jeon and Lee, 2023). As illustrated in Figure 2, this paper explores three critical *roles of LLMs* in this context: their function as **data enhancers** (Section 4) to optimize learning materials, their capacity as **task predictors** (Section 5) to tailor educational solutions, and their potential as **agents** (Section 6) that deliver interactive and adaptive language instruction. By examining these roles, we aim to demonstrate how LLMs can address the limitations of traditional English teaching methods while advancing our understanding of intelligent tutoring systems. Additionally, we discuss potential challenges, ethical considerations, and future directions (Appendix C) for integrating LLMs into English Education, offering a technical guideline for researchers and educators to harness their transformative potential. We also describe the paradigm shift of leveraging AI for English Education, starting from the last century, as one of our contributions in Section 3.

## 2 Background

### 2.1 English Education

Traditional English Education methods often emphasize grammar rules, vocabulary memorization, and repetitive practice, supplemented by limited opportunities for real-world application (Watzke, 2003). Such approaches are often constrained by

the availability of skilled teachers, the diversity of learners’ needs, and the lack of personalized feedback (Williams et al., 2004). Recently, many technologies for English Education have been proposed (Alhusaiyan, 2024), focusing on solving specific tasks instead of describing the whole picture of English tutoring. While intelligent language tutoring systems have the potential to create adaptive environments, attention to this field is relatively less compared to other subjects like science (Shao et al., 2025) and mathematics (Ahn et al., 2024). One key reason lies in the inherent complexity of language as an *ill-defined* domain (Schmidt and Strasser, 2022), posing a great challenge in establishing a valid automatic analysis of learner languages due to the vast variability and unpredictability of human language.

## 2.2 Large Language Models for Education

The potential of LLMs in education (Alhafni et al., 2024), particularly in English Education (Gao et al., 2024; Karataş et al., 2024; Cherednichenko et al., 2024), is immense. Benefiting from large-scale pre-training on extensive corpora, LLMs have demonstrated emergent abilities including (1) *in-context learning* (Dong et al., 2022), which allows the model to adapt to new tasks and provide contextually relevant responses based on a few examples provided during the interaction; (2) *instruction following* (Zeng et al., 2024), which enables the model to process and execute complex user instructions with high accuracy; and (3) *reasoning and planning* (Huang et al., 2024b), which allows the model to generate coherent, structured, and context-aware outputs, even for tasks that require multi-step thinking. However, these fundamental capabilities, while impressive, are insufficient to fully meet the unique demands of English Education. Teaching English requires more than generating grammatically correct sentences or providing accurate translations; it demands a nuanced understanding of pedagogy, learner psychology, and cultural context. Maurya et al. (2024) propose an evaluation taxonomy that identifies eight critical dimensions for assessing AI tutors. These dimensions can be broadly categorized into two groups. (1) *Problem-solving abilities* assess the technical capabilities of LLMs to perform tasks relevant to English Education. (2) *Pedagogical alignment abilities* evaluate how well the LLM aligns with effective teaching and learning principles. Pedagogical alignment includes the model’s ability to adapt to the learner’s

proficiency level, provide scaffolded feedback, foster engagement, and maintain motivation. While LLMs can give direct answers, their ability to replicate these nuanced teaching strategies remains a challenge (Wang et al., 2024a).

## 3 Paradigm Shift

The development of AI models for English Education can be broadly categorized into four successive generations as shown in Figure 3: (1) *rule-based models*, (2) *statistical models*, (3) *neural models*, and (4) *large language models*. We leave the detailed description in Appendix B.

**Our position.** We foresee next-generation LLMs with deeper alignment to pedagogical principles and stronger guardrails to mitigate misinformation and bias. Future models may integrate multimodal data (e.g., text, image, video, speech) to adapt to diverse learner profiles in real time. These improvements will reinforce the position that LLMs can evolve into more effective tutors for English Education.

## 4 LLMs as Data Enhancers

Education is a high-stake area where any hallucination could cause devastating harm to humans’ cognition activities (Ho et al., 2024). One of the hallucination causes is from data (Huang et al., 2023). Therefore, high-quality and diverse data resources (Long et al., 2024) are critical to ensuring the reliability of incorporating LLMs into English Education. The 1) *creation*, 2) *reformation*, and 3) *annotation* of educational materials are crucial to delivering effective and engaging teaching. Traditional resource development methods often lack the scalability, adaptability, and personalization necessary to meet the diverse needs of learners (Feng et al., 2021; Shorten et al., 2021). In contrast, LLMs emerge as transformative tools capable of enhancing these processes (Wang et al., 2024c; Liu et al., 2024c). This section explores how LLMs serve as data enhancers in English Education.

### 4.1 Data Creation

Creating pedagogically sound and learner-specific data is a cornerstone of personalized learning. However, manually creating such resources is time-consuming and often fails to address the wide range of learner needs (Cochran et al., 2022). LLMs can revolutionize this process by generating tailored

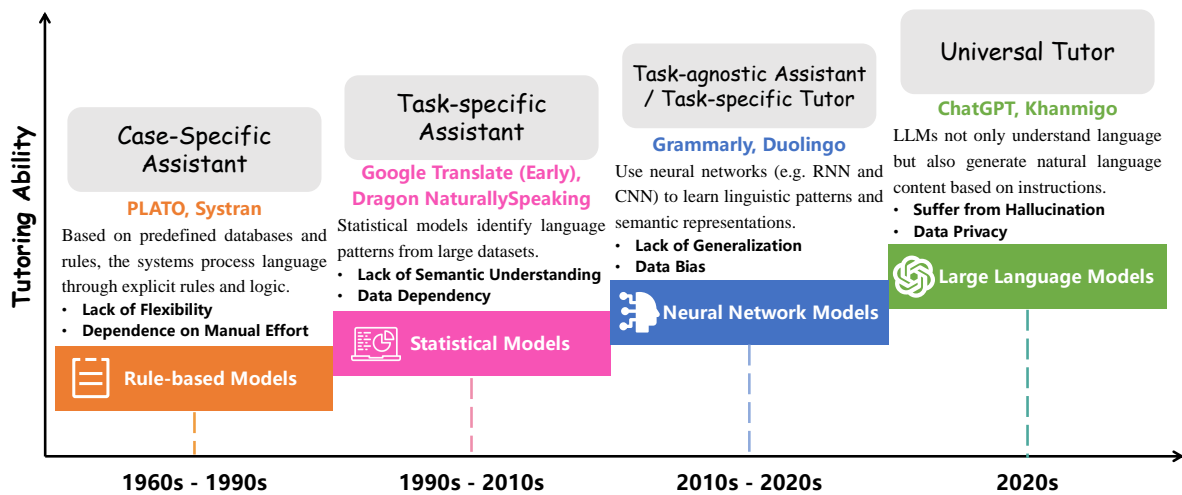


Figure 3: Roadmap of English Education.

and diverse educational content or responses on demand (Zha et al., 2023; Cochran et al., 2023).

**Educational Materials Generation.** A primary use of LLMs in data creation is the *generation of educational questions* aligned with specific learning objectives. Due to their superior contextual understanding, classic rule-based approaches have largely been eclipsed by neural network-based techniques (Kurdi et al., 2020; Rathod et al., 2022; Mulla and Gharpure, 2023). LLMs can produce answer-aware (whose target answer is known) or answer-agnostic (whose answer is open) (Zhang et al., 2021), resulting in more nuanced exercises and assessments (Xiao et al., 2023).

**Student Simulation.** Simulating the learner’s perspective is crucial for designing adaptive instructional materials. Traditional surveys and standardized tests often fail to capture the complexity of dynamic learner behaviors (Käser and Alexandron, 2024). In contrast, LLM-based approaches enable high-fidelity, context-aware *student simulations* (Liu et al., 2024d; Yue et al., 2024), generating synthetic learners who exhibit realistic mastery levels and evolving behaviors. For instance, *Generative Students* (Lu and Wang, 2024) create simulated learners with various competency levels, while *EduAgent* (Xu et al., 2024) integrates cognitive priors to model complex learning trajectories and behaviors better.

**Discussion.** While LLMs excel at generating educational content, current approaches mainly focus on question creation, leaving many areas of English Education underexplored. Essential tasks like generating culturally rich reading materials, context-

dependent writing prompts, or dynamic comprehension exercises still lack diversity and depth. Additionally, the student simulations created by LLMs often fail to reflect long-term learning trajectories or the intricacies of individual learning progress.

## 4.2 Data Reformation

In addition to creating new content, LLMs can adapt *existing* materials to better align with current needs. This process, commonly referred to as data reformation, involves (1) changing data types or modalities, (2) paraphrasing materials to match learner proficiency, and (3) enriching raw data with auxiliary signals or contextual content.

**Teaching Material Transformation.** Transforming existing materials into different forms can yield more comprehensive and immersive learning experiences. For example, *Book2Dial* (Wang et al., 2024b) generates teacher-student dialogues grounded in textbooks, keeping the content both relevant and informative. Their approach includes multi-turn question generation and answering (Kim et al., 2022), dialogue inpainting (Dai et al., 2022), and role-playing. Likewise, *Slide2Lecture* (Zhang-Li et al., 2024) automatically converts lecture slides into structured teaching agendas, enabling interactive follow-up and deeper learner engagement.

**Simplification and Paraphrasing.** Another vital application is simplifying or paraphrasing complex texts to specified readability levels (Huang et al., 2024a) without losing key concepts (Al-Thanyyan and Azmi, 2021). This is particularly beneficial in English Education settings, where language beginners often face advanced vocabulary and complex



structures (Day et al., 2025). Recent advancements in controllable generation (Zhang et al., 2023) leverage model fine-tuning on curated datasets (Zeng et al., 2023) or decoding-time interventions (Liang et al., 2024), thereby allowing educators to specify text complexity, style, or tone.

**Cultural Context Adaptation.** Beyond linguistic correctness, cultural nuance is another crucial factor in English Education (Byram, 1989, 2008). LLMs can facilitate this process by recontextualizing existing materials to reflect the cultural and social norms of different areas (Liu et al., 2024a; Adilazuarda et al., 2024; Kharchenko et al., 2024). For instance, a short story originally set in an English-speaking environment may be adapted for Japanese students by adjusting the characters' names, idiomatic expressions, or social customs, while preserving core instructional goals. This cultural adaptation not only enhances learner engagement but also strengthens cross-cultural competencies.

**Discussion.** While LLM-based data reformation can significantly enhance English Education, several gaps warrant attention. Most current studies prioritize textual forms or single-modal approaches, which may overlook valuable *multimodal* resources such as interactive video and audio-based content (Ghosal et al., 2023). Furthermore, cultural adaptation, although promising, remains underexplored in practical classroom scenarios, particularly for underrepresented persons and culturally sensitive topics. AlKhamissi et al. (2024) demonstrate how cultural misalignment can increase bias. However, robust empirical *evaluations* are still limited across diverse learners and linguistic backgrounds.

### 4.3 Data Annotation

While *Data Creation* focuses on generating learner-specific data, it often prioritizes diversity and adaptability over precision. The approach is particularly useful for tasks with large label spaces (Ding et al., 2024). In contrast, *Data Annotation* emphasizes producing high-quality, meticulously labeled data that is essential for tasks requiring accuracy and consistency. Unlike data creation, annotated data often undergoes rigorous validation to ensure its accuracy and relevancy (Artemova et al., 2024).

**Annotation Generation.** LLMs can be central to generating a variety of annotations, including categorical labels, rationales, pedagogical feedback, and linguistic features such as discourse relations.

Recent prompt engineering and fine-tuning techniques have further expanded LLMs' annotation capabilities. For instance, Ye et al. (2024) leverage GPT-4 to annotate structured explanations for Chinese grammatical error correction, while Samuel et al. (2024) examine GPT-4 as a surrogate for human annotators in low-resource reading comprehension tasks. Likewise, Siyan et al. (2024) deploy GPT-4-Turbo for audio transcript annotations. However, inconsistencies across LLMs (Törnberg, 2024) remain a serious challenge, posing risks to educational reliability.

**Annotation Assessment.** Although LLM-based annotation is efficient, it also raises critical issues of bias, calibration, and validity, particularly in low-resource language contexts (Bhat and Varma, 2023; Jadhav et al., 2024). Automated or semi-automated evaluation strategies have emerged to address these quality concerns. For example, LLMs-as-Judges (Li et al., 2024a,b; Gu et al., 2024) reduce human overhead by automating evaluation, an approach increasingly explored in education-focused applications (Chiang et al., 2024; Zhou et al., 2024). However, purely automated frameworks can still propagate errors or bias.

**Discussion.** Although LLMs provide efficient data annotation, the inconsistency across different models remains a critical concern, affecting the quality and reliability of annotated educational materials. These discrepancies hinder the creation of universally reliable educational content, especially in diverse linguistic and cultural contexts. Additionally, automated annotations often lack the nuance needed for pedagogical applications, making it essential to involve human oversight in critical cases to mitigate errors or biases.

**Our position.** We acknowledge the current limitations in LLM-based data creation, reformation, and annotation for English Education. However, we believe that with continued interdisciplinary collaboration, these challenges can be addressed. *Future advancements* should focus on enhancing the accuracy and diversity of generated content, improving multi-modal and culturally sensitive learning materials, and integrating more robust systems for human-LLM collaboration (Li et al., 2023; Wang et al., 2024e) in data annotation. This will ensure that LLMs can fully realize their potential as effective tutors in English Education.

## 5 LLMs as Task Predictors

*Task-Based Language Learning (TBLL)* (Nunan, 1989; Willis, 2021) as a methodological approach is one of the effective English Education methods. LLMs have demonstrated remarkable capabilities in understanding and generating human language, making them well-suited for addressing numerous tasks in English Education. These tasks can be broadly categorized into three types based on their nature and the role of LLMs: 1) *Discriminative*, 2) *Generative*, and 3) *Mixed* of the above two roles.

### 5.1 Discriminative Task Predictors

Discriminative tasks in English Education primarily involve classifying learner inputs or grading their future performance. Below are some applications that are still calling for improvements:

**Automated Assessment.** The task aims to automatically grade students' assignments, including essay scoring (Seßler et al., 2024; Li and Liu, 2024; Syamkumar et al., 2024), short answer grading (Schneider et al., 2023; Henkel et al., 2024), and spoken language evaluation (Gao et al., 2023; Fu et al., 2024). LLMs can process learners' submissions to judge grammar, lexical diversity, coherence, and even spoken fluency, providing instant feedback. This scalability is particularly appealing for large classes, where human evaluators are often overwhelmed and unable to provide timely, personalized critique (Mizumoto and Eguchi, 2023).

**Knowledge Tracing.** Given sequences of learning interactions in online learning systems, Knowledge Tracing identifies and tracks students' evolving mastery of target skills (Shen et al., 2024b; Xu et al., 2023). LLM-based methods of Knowledge Tracing have been explored in cold-start scenarios (Zhan et al., 2024; Jung et al., 2024), offering strong generalization by inferring latent learner states from limited data. These approaches can support adaptive learning pathways, giving personalized recommendations based on predicted performance and knowledge gaps.

**Discussion.** Despite their promise in automating and personalizing these discriminative tasks, LLMs still grapple with notable limitations that hinder their utility as robust tutoring tools. First, *misalignment of assessment with expert instructors* poses risks: machine-generated scores may deviate from established rubrics or neglect qualitative nuances,

leading to potential discrepancies in grading quality (Kundu and Barbosa, 2024). Second, the *lack of empathy* compounds this issue, as assessments devoid of human judgment risk discouraging learners or overlooking subtle motivational factors (Sharma et al., 2024). Knowledge tracing approaches, while promising in cold-start scenarios, struggle with capturing the complexity of long-term learning trajectories and deeper cognitive processes (Cho et al., 2024). These concerns point to the need for more transparent and human-centered methods in utilizing LLMs for assessment.

### 5.2 Generative Task Predictors

Generative tasks involve producing new content or responses. LLMs are known to be adept at these tasks due to their generation capabilities.

**Grammatical Error Correction and Explanation.** In English writing, errors often reveal learners' gaps in grammar and vocabulary (Hyland and Hyland, 2006). LLMs can detect and correct these errors (Bryant et al., 2023; Ye et al., 2023), offering concise explanations (Ye et al., 2024) that reinforce language rules. By streamlining error detection and corrections, learners deepen their linguistic understanding.

**Feedback Generation.** Quizzes and exercises remain vital in English Education for practice and targeted remediation (Rashov, 2024). LLMs enhance this process by delivering prompt, personalized feedback that pinpoints strengths and addresses weaknesses (Borges et al., 2024). This scalability enables learners to self-regulate and refine their skills without relying solely on human graders (Stamper et al., 2024).

**Socratic Dialogue.** Moving beyond straightforward Q&A, Socratic questioning promotes critical thinking and self-reflection (Paul and Elder, 2007). *SocraticLM* (Liu et al., 2024b), for example, aligns an LLM with open-ended, inquiry-based teaching principles, guiding learners through iterative exploration rather than prescriptive correction. In theory, this fosters deeper conceptual understanding and active learner engagement.

**Discussion.** Despite the promise of LLM-based generation in English Education, multiple uncertainties persist. *Determining how to provide automatic feedback that genuinely maximizes learning outcomes* is an ongoing challenge (Stamper et al., 2024), particularly given education's risk-averse

culture and high accountability standards (Xiao et al., 2024). Moreover, while LLMs like SocraticLM have demonstrated success in domains like mathematics, their applicability to English Education contexts has not been thoroughly validated (Liu et al., 2024b). As such, the design of strategies and follow-up queries remains an open question in ensuring these systems track and respond to learners' cognitive states.

### 5.3 Mixed Task Predictors

Mixed tasks integrate discriminative and generative elements, requiring LLMs to evaluate learner inputs and generate meaningful feedback or suggestions. These tasks are particularly valuable in fostering an interactive and adaptive learning experience, as they bridge the gap between evaluation and instruction.

**Automated Assessment with Feedback.** While discriminative systems for automated essay scoring and speech evaluation primarily focus on assigning grades, LLMs extend these capabilities by simultaneously generating formative feedback (Katuka et al., 2024; Stahl et al., 2024b). For example, an LLM can evaluate the coherence and lexical diversity of a written essay, then offer specific revision strategies. In speaking practice, it can measure fluency and pronunciation accuracy while suggesting drills to refine intonation or stress patterns. Through this combination of scoring and tailored advice, learners gain a deeper understanding of their strengths and areas for improvement.

**Error Analysis.** Error Analysis systematically uncovers and categorizes learners' missteps, from syntactic lapses in writing to flawed pronunciations in speaking (James, 2013; Erdoğan, 2005). LLMs functioning in a mixed capacity can classify these errors and generate corrective guidance, providing revised sentences, clarifications of grammatical rules, or remediation exercises for identified weaknesses (Myles, 2002; Mashoor and Abdulah, 2020). Such insight facilitates targeted interventions that enhance language proficiency across modalities, including reading and listening.

**Discussion.** Mixed-task systems hold promise by combining assessment and feedback generation, but they face notable challenges. One major issue is the *weak alignment* between scoring mechanisms and the quality of feedback provided (Stahl et al., 2024b). For example, while essay scoring

systems may deliver comprehensive evaluations, the feedback often lacks specificity, limiting its instructional value. Additionally, although error analysis has potential, *the absence of standardized pedagogical benchmarks*, especially in oral tasks, hampers the reliability and comparability of LLM-based tools (Leu Jr, 1982).

**Our position.** While LLMs offer scalable solutions for task prediction in English Education, their current limitations—such as misalignment with expert assessments, lack of empathy, and weak alignment between assessment and feedback—require ongoing refinement. *Future research* should focus on improving model transparency, enhancing the cultural and emotional sensitivity of LLMs, and refining task predictors to better reflect long-term learning trajectories and learner motivation. Additionally, developing standardized pedagogical benchmarks for error analysis will help ensure the consistency and reliability of LLM-generated feedback.

## 6 LLM-empowered Agent

In this section, we delve into the potential of LLMs as intelligent tutoring agents in English Education. LLMs can act as catalysts for personalized learning, addressing the long-standing scalability, adaptability, and inclusivity challenges in traditional teaching paradigms.

### 6.1 Fundamental Abilities

This section highlights five key abilities of LLM-empowered agents that enable them to function as adaptive tutors.

**Knowledge Integration.** LLMs excel at merging structured educational knowledge graphs (Abu-Rasheed et al., 2024; Hu and Wang, 2024) with unstructured textual data (Li et al., 2024c; Modran et al., 2024), providing rich, contextualized information on linguistic constructs and cultural nuances. Their ability to perform real-time knowledge editing (Wang et al., 2024d; Zhang et al., 2024a) ensures learners receive content aligned with evolving language usage, addressing the inherent limitations of static materials.

**Pedagogical Alignment.** LLMs require embedding with pedagogical principles to facilitate genuine learning experiences (Carroll, 1965; Taneja, 1995). Recent work incorporates theoretical frameworks, such as Bloom's taxonomy (Bloom et al.,



1956), to guide LLMs in systematically addressing different cognitive levels (Jiang et al., 2024b). Approaches like *Pedagogical Chain of Thought* (Jiang et al., 2024b) and *preference learning* (Sonkar et al., 2024; Rafailov et al., 2024) focus on aligning model responses with educational objectives.

**Planning.** By assisting in crafting teaching objectives and lesson designs, LLMs can handle complex tasks such as differentiated instruction (Hu et al., 2024). LessonPlanner (Fan et al., 2024) has been proposed to assist novice teachers in preparing lesson plans, with expert interviews confirming its effectiveness. Zheng et al. (2024) propose a three-stage process to produce customized lesson plans, using Retrieval-Augmented Generation (RAG), self-critique, and subsequent refinement.

**Memory.** Effective tutoring systems track learner histories and tailor subsequent interactions accordingly (Jiang et al., 2024a; Chen et al., 2024). When serving as memory-augmented agents, LLMs can retain individualized data—such as repeated grammar mistakes or overlooked vocabulary—thereby improving continuity and enabling consistent scaffolding of future learning tasks.

**Tool Using.** Beyond textual interactions, LLM-based agents can integrate specialized tools to streamline the educational ecosystem, from cognitive diagnosis modules (Ma and Guo, 2019) to report generators (Zhou et al., 2025). By orchestrating these resources, LLMs seamlessly unify diverse utilities under a single interface, enhancing learner experience and instructional efficiency.

## 6.2 Applications

Although still in its early stages, LLM-empowered agents have already started to show promising applications in English Education.

**Classroom Simulation.** Classroom simulation leverages LLM-empowered agents to recreate complex, interactive learning settings without the logistical hurdles of organizing physical classrooms (Zhang et al., 2024b). By simulating virtual students and tutors, researchers can study pedagogical strategies at scale, generate diverse learner interactions, and refine teaching techniques. Moreover, this virtual data can be used to fine-tune LLMs for specific educational contexts and learner profiles (Liu et al., 2024b), offering a cost-effective and adaptable approach to language instruction.

**Intelligent Tutoring System (ITS).** LLM-based agents have demonstrated the capacity to provide dynamic, personalized tutoring experiences (Wang et al., 2025; Kwon et al., 2024), effectively identifying learner weaknesses through large-scale linguistic analysis (Caines et al., 2023). This makes them promising for delivering individualized instruction at scale. Although current ITS applications in mathematics (Pal Chowdhury et al., 2024) and science (Stamper et al., 2024) have shown success, the extension to English Education requires nuanced handling of cultural and contextual elements, as well as the unpredictability of human language usage.

**Discussion.** Despite the promise of these applications, critical challenges remain. Existing classroom simulation frameworks often *lack standardized benchmarks for English Education*, making it difficult to assess the efficacy and generalizability of developed systems (Zhang et al., 2024b). In addition, evaluating language-specific tutoring strategies, including real-time conversational practice and holistic skill integration, remains an underexplored frontier. Addressing these gaps requires *new datasets and metrics* centered on holistic skill development and interdisciplinary collaboration.

**Our position.** We argue that *future research* should focus on integrating multimodal learning tasks (Sonlu et al., 2024) and developing standardized frameworks for evaluating English Education simulations. Moreover, LLMs should evolve beyond text-based capabilities to provide real-time, context-sensitive feedback, particularly in speaking and listening. Interdisciplinary collaboration and the creation of new datasets tailored to English Education are crucial for refining these systems and ensuring their scalability and inclusivity in language instruction.

## 7 Conclusion

This paper emphasizes the transformative potential of LLMs in English Education, positioning them as valuable tutors to complement traditional teaching methods. Through their roles as data enhancers, task predictors, and agents, LLMs can provide adaptive learning experiences across the core skills of listening, speaking, reading, and writing. This paper encourage continuing dialogue and interdisciplinary collaboration to responsibly integrate LLMs into educational ecosystems.



## Limitations

**Emphasis on potential over practical implementation barriers.** This paper primarily focuses on the potential of LLMs to serve as effective tutors in English Education, outlining beneficial roles as data enhancers, task predictors, and agents. While we acknowledge the existence of challenges (to be discussed in Appendix C), a limitation of this position is that the main arguments may not fully capture the considerable practical, socio-economic, and infrastructural hurdles that could impede the equitable and effective implementation of these LLM roles across diverse global educational contexts and resource settings.

**Generalizability and contextual adaptation of proposed roles.** We propose three broad roles for LLMs in English Education. However, this paper does not provide an exhaustive analysis of how the efficacy and suitability of LLMs in these roles might vary significantly across different target languages (especially low-resource languages), specific learner demographics (e.g., preschoolers vs. K-12 vs. adult learners, learners with disabilities), diverse cultural contexts, or varying pedagogical philosophies. The general framework presented may require substantial adaptation and further research to be effectively applied in specific English Education scenarios.

**Nuances of human-LLM pedagogical interaction.** While advocating for LLMs as tutors that can complement human expertise, this position paper does not delve deeply into the complex dynamics of the pedagogical interactions between learners, LLM-based tutors, human educators, and parents. Critical aspects such as optimizing the collaborative model, designing effective training for educators to leverage LLMs, mitigating risks of learner over-reliance, and ensuring that LLM interactions foster deep learning rather than superficial engagement are multifaceted issues that warrant more extensive investigation than afforded by the scope of this paper.

## References

Hasan Abu-Rasheed, Christian Weber, and Madjid Fathi. 2024. Knowledge graphs as context sources for llm-based explanations of learning recommendations. *arXiv preprint arXiv:2403.03008*.

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,

Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Muhammad Farid Adilazuarda, Sagnik Mukherjee, Pradhyumna Lavania, Siddhant Shivdutt Singh, Alham Fikri Aji, Jacki O’Neill, Ashutosh Modi, and Monojit Choudhury. 2024. *Towards measuring and modeling “culture” in LLMs: A survey*. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 15763–15784, Miami, Florida, USA. Association for Computational Linguistics.

Janice Ahn, Rishu Verma, Renze Lou, Di Liu, Rui Zhang, and Wenpeng Yin. 2024. Large language models for mathematical reasoning: Progresses and challenges. *arXiv preprint arXiv:2402.00157*.

Suha S Al-Thanyyan and Aqil M Azmi. 2021. Automated text simplification: a survey. *ACM Computing Surveys (CSUR)*, 54(2):1–36.

Bashar Alhafni, Sowmya Vajjala, Stefano Bannò, Kaushal Kumar Maurya, and Ekaterina Kochmar. 2024. Llms in education: Novel perspectives, challenges, and opportunities. *arXiv preprint arXiv:2409.11917*.

Eman Alhusaiyan. 2024. A systematic review of current trends in artificial intelligence in foreign language learning. *Saudi Journal of Language Studies*.

Eman Alhusaiyan. 2025. A systematic review of current trends in artificial intelligence in foreign language learning. *Saudi Journal of Language Studies*, 5(1):1–16.

Badr AlKhamissi, Muhammad ElNokrashy, Mai Alkhamissi, and Mona Diab. 2024. *Investigating cultural alignment of large language models*. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12404–12422, Bangkok, Thailand. Association for Computational Linguistics.

Preeti Anand. 2023. Khan academy creates gpt-4 based helper khanmigo marking formal entry of ai into education.

Ekaterina Artemova, Akim Tsvigun, Dominik Schlechtweg, Natalia Fedorova, Sergei Tilga, and Boris Obmoroshev. 2024. Hands-on tutorial: Labeling with llm and human-in-the-loop. *arXiv preprint arXiv:2411.04637*.

Sumit Asthana, Hannah Rashkin, Elizabeth Clark, Fantine Huot, and Mirella Lapata. 2024. *Evaluating LLMs for targeted concept simplification for domain-specific texts*. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 6208–6226, Miami, Florida, USA. Association for Computational Linguistics.

739	Savita Bhat and Vasudeva Varma. 2023. <a href="#">Large language models as annotators: A preliminary evaluation for annotating low-resource language content</a> . In <i>Proceedings of the 4th Workshop on Evaluation and Comparison of NLP Systems</i> , pages 100–107, Bali, Indonesia. Association for Computational Linguistics.	794
740		795
741		796
742		797
743		
744		
745		
746	Christopher Blair. 1997. Dragon–naturallyspeaking. <i>Journal of Osteopathic Medicine</i> , 97(12):711–711.	
747		
748	Benjamin S Bloom, Max D Engelhart, Edward J Furst, Walker H Hill, David R Krathwohl, and 1 others. 1956. <i>Taxonomy of educational objectives: The classification of educational goals. Handbook 1: Cognitive domain</i> . Longman New York.	
749		
750		
751		
752		
753	Angana Borah and Rada Mihalcea. 2024. Towards implicit bias detection and mitigation in multi-agent llm interactions. <i>arXiv preprint arXiv:2410.02584</i> .	
754		
755		
756	Beatriz Borges, Niket Tandon, Tanja Käser, and Antoine Bosselut. 2024. <a href="#">Let me teach you: Pedagogical foundations of feedback for language models</a> . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 12082–12104, Miami, Florida, USA. Association for Computational Linguistics.	
757		
758		
759		
760		
761		
762		
763	Christopher Bryant, Zheng Yuan, Muhammad Reza Qorib, Hannan Cao, Hwee Tou Ng, and Ted Briscoe. 2023. Grammatical error correction: A survey of the state of the art. <i>Computational Linguistics</i> , 49(3):643–701.	
764		
765		
766		
767		
768	M Byram. 1989. Cultural studies in foreign language education. <i>Multilingual Matters</i> , 61.	
769		
770	Michael Byram. 2008. <i>From foreign language education to education for intercultural citizenship: Essays and reflections</i> , volume 17. Multilingual matters.	
771		
772		
773	Marios C Angelides and Isabel Garcia. 1993. Towards an intelligent knowledge based tutoring system for foreign language learning. <i>Journal of computing and information technology</i> , 1(1):15–28.	
774		
775		
776		
777	Andrew Caines, Luca Benedetto, Shiva Taslimipoor, Christopher Davis, Yuan Gao, Oeistein Andersen, Zheng Yuan, Mark Elliott, Russell Moore, Christopher Bryant, and 1 others. 2023. On the application of large language models for language teaching and assessment technology. <i>arXiv preprint arXiv:2307.08393</i> .	
778		
779		
780		
781		
782		
783		
784	John B Carroll. 1965. The contributions of psychological theory and educational research to the teaching of foreign languages. <i>The modern language journal</i> , 49(5):273–281.	
785		
786		
787		
788	Yulin Chen, Ning Ding, Hai-Tao Zheng, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. 2024. Empowering private tutoring by chaining large language models. In <i>Proceedings of the 33rd ACM International Conference on Information and Knowledge Management</i> , pages 354–364.	
789		
790		
791		
792		
793		
	Olga Cherednichenko, Olha Yanholenko, Antonina Badan, Nataliia Onishchenko, and Nunu Akopiants. 2024. Large language models for foreign language acquisition.	794
		795
		796
		797
	Cheng-Han Chiang, Wei-Chih Chen, Chun-Yi Kuan, Chienchou Yang, and Hung-yi Lee. 2024. <a href="#">Large language model as an assignment evaluator: Insights, feedback, and challenges in a 1000+ student course</a> . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 2489–2513, Miami, Florida, USA. Association for Computational Linguistics.	798
		799
		800
		801
		802
		803
		804
		805
	Yongwan Cho, Rabia Emhamed AlMamlook, and Tasnim Gharaibeh. 2024. A systematic review of knowledge tracing and large language models in education: Opportunities, issues, and future research. <i>arXiv preprint arXiv:2412.09248</i> .	806
		807
		808
		809
		810
	Keith Cochran, Clayton Cohn, Nicole Hutchins, Gautam Biswas, and Peter Hastings. 2022. Improving automated evaluation of formative assessments with text data augmentation. In <i>International Conference on Artificial Intelligence in Education</i> , pages 390–401. Springer.	811
		812
		813
		814
		815
		816
	Keith Cochran, Clayton Cohn, Jean Francois Rouet, and Peter Hastings. 2023. Improving automated evaluation of student text responses using gpt-3.5 for text data augmentation. In <i>International Conference on Artificial Intelligence in Education</i> , pages 217–228. Springer.	817
		818
		819
		820
		821
		822
	Zhuyun Dai, Arun Tejasvi Chaganty, Vincent Y Zhao, Aida Amini, Qazi Mamunur Rashid, Mike Green, and Kelvin Guu. 2022. Dialog inpainting: Turning documents into dialogs. In <i>International conference on machine learning</i> , pages 4558–4586. PMLR.	823
		824
		825
		826
		827
	Stephanie L. Day, Jacapo Cirica, Steven R. Clapp, Veronika Penkova, Amy E. Giroux, Abbey Banta, Catherine Bordeau, Poojitha Muteneni, and Ben D. Sawyer. 2025. <a href="#">Evaluating genai for simplifying texts for education: Improving accuracy and consistency for enhanced readability</a> . <i>Preprint</i> , arXiv:2501.09158.	828
		829
		830
		831
		832
		833
		834
	Bosheng Ding, Chengwei Qin, Ruochen Zhao, Tianze Luo, Xinze Li, Guizhen Chen, Wenhan Xia, Junjie Hu, Anh Tuan Luu, and Shafiq Joty. 2024. <a href="#">Data augmentation using LLMs: Data perspectives, learning paradigms and challenges</a> . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 1679–1705, Bangkok, Thailand. Association for Computational Linguistics.	835
		836
		837
		838
		839
		840
		841
		842
	Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Tianyu Liu, and 1 others. 2022. A survey on in-context learning. <i>arXiv preprint arXiv:2301.00234</i> .	843
		844
		845
		846
	Sarah Elaine Eaton. 2010. <i>Global Trends in Language Learning in the 21st Century</i> . ERIC.	847
		848

849	Ronald G Ehrenberg, Dominic J Brewer, Adam	Yingming Gao, Baorian Nuchged, Ya Li, and Linkai	902
850	Gamoran, and J Douglas Willms. 2001. Class size	Peng. 2023. An investigation of applying large lan-	903
851	and student achievement. <i>Psychological science in</i>	guage models to spoken language learning. <i>Applied</i>	904
852	<i>the public interest</i> , 2(1):1–30.	<i>Sciences</i> , 14(1):224.	905
853	Vacide Erdoğan. 2005. Contribution of error analysis	Deepanway Ghosal, Navonil Majumder, Ambuj	906
854	to foreign language teaching. <i>Mersin Üniversitesi</i>	Mehrish, and Soujanya Poria. 2023. Text-to-audio	907
855	<i>Eğitim Fakültesi Dergisi</i> , 1(2).	generation using instruction-tuned llm and latent dif-	908
856		fusion model. <i>arXiv preprint arXiv:2304.13731</i> .	909
857	Haoxiang Fan, Guanzheng Chen, Xingbo Wang, and	Crina Grosan, Ajith Abraham, Crina Grosan, and Ajith	910
858	Zhenhui Peng. 2024. Lessonplanner: Assisting	Abraham. 2011. Rule-based expert systems. <i>Intelli-</i>	911
859	novice teachers to prepare pedagogy-driven lesson	<i>gent systems: A modern approach</i> , pages 149–185.	912
860	plans with large language models. In <i>Proceedings of</i>		
861	<i>the 37th Annual ACM Symposium on User Interface</i>	Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan,	913
	<i>Software and Technology</i> , pages 1–20.	Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan	914
862	Tao Fang, Shu Yang, Kaixin Lan, Derek F Wong, Jin-	Shen, Shengjie Ma, Honghao Liu, and 1 others.	915
863	peng Hu, Lidia S Chao, and Yue Zhang. 2023. Is	2024. A survey on llm-as-a-judge. <i>arXiv preprint</i>	916
864	chatgpt a highly fluent grammatical error correction	<i>arXiv:2411.15594</i> .	917
865	system? a comprehensive evaluation. <i>arXiv preprint</i>		
866	<i>arXiv:2304.01746</i> .	Jieun Han, Haneul Yoo, Yoonsu Kim, Junho Myung,	918
867	Lucile Favero, Juan Antonio Pérez-Ortiz, Tanja Käser,	Minsun Kim, Hyunseung Lim, Juho Kim, Tak Yeon	919
868	and Nuria Oliver. 2024. Enhancing critical thinking	Lee, Hwajung Hong, So-Yeon Ahn, and 1 others.	920
869	in education by means of a socratic chatbot. <i>arXiv</i>	2023a. Recipe: How to integrate chatgpt into efl	921
870	<i>preprint arXiv:2409.05511</i> .	writing education. In <i>Proceedings of the tenth ACM</i>	922
		<i>conference on learning@ scale</i> , pages 416–420.	923
871	Yuejiao Fei, Leyang Cui, Sen Yang, Wai Lam, Zhen-	Jieun Han, Haneul Yoo, Junho Myung, Minsun Kim,	924
872	zhong Lan, and Shuming Shi. 2023. <a href="#">Enhancing gram-</a>	Hyunseung Lim, Yoonsu Kim, Tak Yeon Lee, Hwa-	925
873	<a href="#">matical error correction systems with explanations.</a>	jung Hong, Juho Kim, So-Yeon Ahn, and Alice Oh.	926
874	In <i>Proceedings of the 61st Annual Meeting of the</i>	2024. <a href="#">LLM-as-a-tutor in EFL writing education: Fo-</a>	927
875	<i>Association for Computational Linguistics (Volume</i>	<a href="#">cusing on evaluation of student-LLM interaction.</a> In	928
876	<i>1: Long Papers</i> ), pages 7489–7501, Toronto, Canada.	<i>Proceedings of the 1st Workshop on Customizable</i>	929
877	Association for Computational Linguistics.	<i>NLP: Progress and Challenges in Customizing NLP</i>	930
878	Steven Y. Feng, Varun Gangal, Jason Wei, Sarath Chan-	<i>for a Domain, Application, Group, or Individual</i>	931
879	dar, Soroush Vosoughi, Teruko Mitamura, and Ed-	<i>(CustomNLP4U)</i> , pages 284–293, Miami, Florida,	932
880	uard Hovy. 2021. <a href="#">A survey of data augmentation</a>	USA. Association for Computational Linguistics.	933
881	<a href="#">approaches for NLP.</a> In <i>Findings of the Association</i>	Jieun Han, Haneul Yoo, Junho Myung, Minsun Kim,	934
882	<i>for Computational Linguistics: ACL-IJCNLP 2021</i> ,	Hyunseung Lim, Yoonsu Kim, Tak Yeon Lee, Hwa-	935
883	pages 968–988, Online. Association for Computa-	jung Hong, Juho Kim, So-Yeon Ahn, and 1 others.	936
884	tional Linguistics.	2023b. Fabric: Automated scoring and feedback gen-	937
885	Tira Nur Fitria. 2021. Grammarly as ai-powered english	eration for essays. <i>arXiv preprint arXiv:2310.05191</i> .	938
886	writing assistant: Students’ alternative for writing		
887	english. <i>Metathesis: Journal of English Language,</i>	Robert Hart. 1981. Language study and the plato sys-	939
888	<i>Literature, and Teaching</i> , 5(1):65–78.	tem. <i>Studies in language learning</i> , 3(1):1–24.	940
889	Nils Freyer, Hendrik Kempt, and Lars Klöser. 2024.	Lin He and Keqin Li. 2024. Mitigating hallucinations	941
890	Easy-read and large language models: on the ethical	in llm using k-means clustering of synonym semantic	942
891	dimensions of llm-based text simplification. <i>Ethics</i>	relevance. <i>Authorea Preprints</i> .	943
892	<i>and Information Technology</i> , 26(3):50.		
893	Kaiqi Fu, Linkai Peng, Nan Yang, and Shuran	Owen Henkel, Libby Hills, Adam Boxer, Bill Roberts,	944
894	Zhou. 2024. Pronunciation assessment with multi-	and Zach Levonian. 2024. Can large language mod-	945
895	modal large language models. <i>arXiv preprint</i>	els make the grade? an empirical study evaluating	946
896	<i>arXiv:2407.09209</i> .	llms ability to mark short answer questions in k-12	947
897	Yang Gao, Qikai Wang, and Xiaochen Wang. 2024. Ex-	education. In <i>Proceedings of the Eleventh ACM Con-</i>	948
898	ploring efl university teachers’ beliefs in integrating	<i>ference on Learning@ Scale</i> , pages 300–304.	949
899	chatgpt and other large language models in language		
900	education: a study in china. <i>Asia Pacific Journal of</i>	Huu-Tuong Ho, Duc-Tin Ly, and Luong Vuong Nguyen.	950
901	<i>Education</i> , 44(1):29–44.	2024. Mitigating hallucinations in large language	951
		models for educational application. In <i>2024 IEEE</i>	952
		<i>International Conference on Consumer Electronics-</i>	953
		<i>Asia (ICCE-Asia)</i> , pages 1–4. IEEE.	954





1065	<i>Chapter of the Association for Computational Linguistics: System Demonstrations</i> , pages 168–176, St. Julians, Malta. Association for Computational Linguistics.	1122
1066		1123
1067		1124
1068		1125
1069	HYEJI KIM, Jongyoul Park, Hyeongbae Jeon, Sidney S Fels, Samuel Dodson, and Kyoungwon Seo. 2025. Augmented educators and ai: Shaping the future of human-ai collaboration in learning. In <i>Proceedings of the Extended Abstracts of the CHI Conference on Human Factors in Computing Systems</i> , pages 1–6.	1126
1070		1127
1071		
1072		1128
1073		1129
1074		1130
1075	Anindita Kundu and Denilson Barbosa. 2024. Are large language models good essay graders? <i>arXiv preprint arXiv:2409.13120</i> .	1131
1076		1132
1077		1133
1078		1134
1079	Ghader Kurdi, Jared Leo, Bijan Parsia, Uli Sattler, and Salam Al-Emari. 2020. A systematic review of automatic question generation for educational purposes. <i>International Journal of Artificial Intelligence in Education</i> , 30:121–204.	1135
1080		
1081		1136
1082		1137
1083	Soonwoo Kwon, Sojung Kim, Minju Park, Seunghyun Lee, and Kyuseok Kim. 2024. <i>BIPED: Pedagogically informed tutoring system for ESL education</i> . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3389–3414, Bangkok, Thailand. Association for Computational Linguistics.	1138
1084		1139
1085		
1086		1140
1087		1141
1088		1142
1089		1143
1090	Jeongmin Lee, Jin-Xia Huang, Minsoo Cho, Yoon-Hyung Roh, Oh-Woog Kwon, and Yunkeun Lee. 2024a. Developing conversational intelligent tutoring for speaking skills in second language learning. In <i>International Conference on Intelligent Tutoring Systems</i> , pages 131–148. Springer.	1144
1091		1145
1092		1146
1093		
1094		1147
1095		1148
1096	Minhwa Lee, Zae Myung Kim, Vivek Khetan, and Dongyeop Kang. 2024b. Human-ai collaborative taxonomy construction: A case study in profession-specific writing assistants. In <i>Proceedings of the Third Workshop on Intelligent and Interactive Writing Assistants</i> , pages 51–57.	1149
1097		1150
1098		1151
1099		1152
1100		1153
1101		1154
1102	Unggi Lee, Haewon Jung, Younghoon Jeon, Younghoon Sohn, Wonhee Hwang, Jewoong Moon, and Hyeoncheol Kim. 2024c. Few-shot is enough: exploring chatgpt prompt engineering method for automatic question generation in english education. <i>Education and Information Technologies</i> , 29(9):11483–11515.	1155
1103		1156
1104		1157
1105		1158
1106		1159
1107		1160
1108	Donald J Leu Jr. 1982. Oral reading error analysis: A critical review of research and application. <i>Reading Research Quarterly</i> , pages 420–437.	1161
1109		1162
1110		1163
1111	Dawei Li, Bohan Jiang, Liangjie Huang, Alimohammad Beigi, Chengshuai Zhao, Zhen Tan, Amrita Bhattacharjee, Yuxuan Jiang, Canyu Chen, Tianhao Wu, and 1 others. 2024a. From generation to judgment: Opportunities and challenges of llm-as-a-judge. <i>arXiv preprint arXiv:2411.16594</i> .	1164
1112		1165
1113		1166
1114		
1115		1167
1116		1168
1117	Haitao Li, Qian Dong, Junjie Chen, Huixue Su, Yujia Zhou, Qingyao Ai, Ziyi Ye, and Yiqun Liu. 2024b. Llm-as-judges: A comprehensive survey on llm-based evaluation methods. <i>arXiv preprint arXiv:2412.05579</i> .	1169
1118		1170
1119		1171
1120		1172
1121		1173
		1174
		1175
		1176
	Kunze Li and Yu Zhang. 2024. <i>Planning first, question second: An LLM-guided method for controllable question generation</i> . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 4715–4729, Bangkok, Thailand. Association for Computational Linguistics.	
	Minzhi Li, Taiwei Shi, Caleb Ziems, Min-Yen Kan, Nancy Chen, Zhengyuan Liu, and Diyi Yang. 2023. <i>CoAnnotating: Uncertainty-guided work allocation between human and large language models for data annotation</i> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 1487–1505, Singapore. Association for Computational Linguistics.	
	Wenchao Li and Haitao Liu. 2024. Applying large language models for automated essay scoring for non-native japanese. <i>Humanities and Social Sciences Communications</i> , 11(1):1–15.	
	Xiu Li, Aron Henriksson, Martin Duneld, Jalal Nouri, and Yongchao Wu. 2024c. Supporting teaching-to-the-curriculum by linking diagnostic tests to curriculum goals: Using textbook content as context for retrieval-augmented generation with large language models. In <i>International Conference on Artificial Intelligence in Education</i> , pages 118–132. Springer.	
	Xun Liang, Hanyu Wang, Yezhaohui Wang, Shichao Song, Jiawei Yang, Simin Niu, Jie Hu, Dan Liu, Shunyu Yao, Feiyu Xiong, and 1 others. 2024. Controllable text generation for large language models: A survey. <i>arXiv preprint arXiv:2408.12599</i> .	
	Chen Cecilia Liu, Iryna Gurevych, and Anna Korhonen. 2024a. Culturally aware and adapted nlp: A taxonomy and a survey of the state of the art. <i>arXiv preprint arXiv:2406.03930</i> .	
	Jiayu Liu, Zhenya Huang, Tong Xiao, Jing Sha, Jinze Wu, Qi Liu, Shijin Wang, and Enhong Chen. 2024b. SocratiClm: Exploring socratic personalized teaching with large language models. In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems</i> .	
	Ruibo Liu, Jerry Wei, Fangyu Liu, Chenglei Si, Yanzhe Zhang, Jinmeng Rao, Steven Zheng, Daiyi Peng, Diyi Yang, Denny Zhou, and 1 others. 2024c. Best practices and lessons learned on synthetic data for language models. <i>arXiv preprint arXiv:2404.07503</i> .	
	Zhengyuan Liu, Stella Xin Yin, Geyu Lin, and Nancy F Chen. 2024d. Personality-aware student simulation for conversational intelligent tutoring systems. <i>arXiv preprint arXiv:2404.06762</i> .	
	Zhexiong Liu, Diane Litman, Elaine Wang, Tianwen Li, Mason Gobat, Lindsay Clare Matsumura, and Richard Correnti. 2025. erewise+ rf: A writing evaluation system for assessing student essay revisions and providing formative feedback. <i>arXiv preprint arXiv:2501.00715</i> .	

1177	Lin Long, Rui Wang, Ruixuan Xiao, Junbo Zhao, Xiao	Seyed Parsa Neshaei, Richard Lee Davis, Adam Haz-	1233
1178	Ding, Gang Chen, and Haobo Wang. 2024. <a href="#">On</a>	imeh, Bojan Lazarevski, Pierre Dillenbourg, and	1234
1179	<a href="#">LLMs-driven synthetic data generation, curation, and</a>	Tanja Käser. 2024. Towards modeling learner perfor-	1235
1180	<a href="#">evaluation: A survey</a> . In <i>Findings of the Association</i>	mance with large language models. <i>arXiv preprint</i>	1236
1181	<i>for Computational Linguistics: ACL 2024</i> , pages	<i>arXiv:2403.14661</i> .	1237
1182	11065–11082, Bangkok, Thailand. Association for		
1183	Computational Linguistics.		
1184	Qingyu Lu, Baopu Qiu, Liang Ding, Kanjian Zhang,	Diane Nicholls, Andrew Caines, and Paula Buttery.	1238
1185	Tom Kocmi, and Dacheng Tao. 2024. <a href="#">Error analysis</a>	2024. The write & improve corpus 2024: Error-	1239
1186	<a href="#">prompting enables human-like translation evaluation</a>	annotated and cefr-labelled essays by learners of en-	1240
1187	<a href="#">in large language models</a> . In <i>Findings of the Asso-</i>	glish.	1241
1188	<i>ciation for Computational Linguistics: ACL 2024</i> ,		
1189	pages 8801–8816, Bangkok, Thailand. Association	David Nunan. 1989. <i>Designing tasks for the commu-</i>	1242
1190	for Computational Linguistics.	<i>nicaive classroom</i> . Cambridge university press.	1243
1191	Xinyi Lu and Xu Wang. 2024. Generative students: Us-	Franz Och. 2006. <a href="#">Statistical machine translation live</a> .	1244
1192	ing llm-simulated student profiles to support question		
1193	item evaluation. <i>arXiv preprint arXiv:2405.11591</i> .	Sankalan Pal Chowdhury, Vilém Zouhar, and Mrinmaya	1245
1194	Wenchao Ma and Wenjing Guo. 2019. Cognitive	Sachan. 2024. Autotutor meets large language mod-	1246
1195	diagnosis models for multiple strategies. <i>British</i>	els: A language model tutor with rich pedagogy and	1247
1196	<i>Journal of Mathematical and Statistical Psychology</i> ,	guardrails. In <i>Proceedings of the Eleventh ACM Con-</i>	1248
1197	72(2):370–392.	<i>ference on Learning@ Scale</i> , pages 5–15.	1249
1198	Jakub Macina, Nico Daheim, Ido Hakimi, Manu Kapur,	Richard Paul and Linda Elder. 2007. Critical thinking:	1250
1199	Iryna Gurevych, and Mrinmaya Sachan. 2025. Math-	The art of socratic questioning. <i>Journal of develop-</i>	1251
1200	tutorbench: A benchmark for measuring open-ended	<i>mental education</i> , 31(1):36.	1252
1201	pedagogical capabilities of llm tutors. <i>arXiv preprint</i>		
1202	<i>arXiv:2502.18940</i> .	Andrew Radford, Martin Atkinson, David Britain, Har-	1253
1203	Bakheet Bayan Nayif Mashoor and ATH bin Abdul-	ald Clahsen, and Andrew Spencer. 2009. <i>Linguistics:</i>	1254
1204	lah. 2020. Error analysis of spoken english language	<i>an introduction</i> . Cambridge University Press.	1255
1205	among jordanian secondary school students. <i>Interna-</i>		
1206	<i>tional Journal of Education and Research</i> , 8(5):75–	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christo-	1256
1207	82.	pher D Manning, Stefano Ermon, and Chelsea Finn.	1257
1208	Kaushal Kumar Maurya, KV Srivatsa, Kseniia	2024. Direct preference optimization: Your language	1258
1209	Petukhova, and Ekaterina Kochmar. 2024. Unify-	model is secretly a reward model. <i>Advances in Neu-</i>	1259
1210	ing ai tutor evaluation: An evaluation taxonomy for	<i>ral Information Processing Systems</i> , 36.	1260
1211	pedagogical ability assessment of llm-powered ai tu-		
1212	tors. <i>arXiv preprint arXiv:2412.09416</i> .	Oybek Rashov. 2024. Modern methods of teaching	1261
1213	Atsushi Mizumoto and Masaki Eguchi. 2023. Exploring	foreign languages. In <i>International Scientific and</i>	1262
1214	the potential of using an ai language model for auto-	<i>Current Research Conferences</i> , pages 158–164.	1263
1215	mated essay scoring. <i>Research Methods in Applied</i>		
1216	<i>Linguistics</i> , 2(2):100050.	Manav Rathod, Tony Tu, and Katherine Stasaski. 2022.	1264
1217	Horia Modran, Ioana Corina Bogdan, Doru Ursuțiu,	<a href="#">Educational multi-question generation for reading</a>	1265
1218	Cornel Samoila, and Paul Livius Modran. 2024. Llm	<a href="#">comprehension</a> . In <i>Proceedings of the 17th Work-</i>	1266
1219	intelligent agent tutoring in higher education courses	<i>shop on Innovative Use of NLP for Building Edu-</i>	1267
1220	using a rag approach. <i>Preprints 2024</i> , 2024070519.	<i>cational Applications (BEA 2022)</i> , pages 216–223,	1268
1221	Nikahat Mulla and Prachi Gharpure. 2023. Auto-	Seattle, Washington. Association for Computational	1269
1222	matic question generation: a review of methodolo-	Linguistics.	1270
1223	gies, datasets, evaluation metrics, and applications.		
1224	<i>Progress in Artificial Intelligence</i> , 12(1):1–32.	Mahefa Abel Razafinirina, William Germain Dimbisoa,	1271
1225	Johanne Myles. 2002. Second language writing and	and Thomas Mahatody. 2024. Pedagogical align-	1272
1226	research: The writing process and error analysis in	ment of large language models (llm) for personalized	1273
1227	student texts. <i>Tesl-ej</i> , 6(2):1–20.	learning: A survey, trends and challenges. <i>Jour-</i>	1274
1228	Mahjabin Nahar, Haeseung Seo, Eun-Ju Lee, Aiping	<i>nal of Intelligent Learning Systems and Applications</i> ,	1275
1229	Xiong, and Dongwon Lee. 2024. <a href="#">Fakes of varying</a>	16(4):448–480.	1276
1230	<a href="#">shades: How warning affects human perception and</a>	Vinay Samuel, Houda Aynaou, Arijit Chowdhury,	1277
1231	<a href="#">engagement regarding LLM hallucinations</a> . In <i>First</i>	Karthik Venkat Ramanan, and Aman Chadha. 2024.	1278
1232	<i>Conference on Language Modeling</i> .	<a href="#">Can LLMs augment low-resource reading compre-</a>	1279
		<a href="#">hension datasets? opportunities and challenges</a> . In	1280
		<i>Proceedings of the 62nd Annual Meeting of the Asso-</i>	1281
		<i>ciation for Computational Linguistics (Volume 4: Stu-</i>	1282
		<i>dent Research Workshop)</i> , pages 307–317, Bangkok,	1283
		Thailand. Association for Computational Linguistics.	1284
		Alexander Scarlatos and Andrew Lan. 2024. Exploring	1285
		knowledge tracing in tutor-student dialogues. <i>arXiv</i>	1286
		<i>preprint arXiv:2409.16490</i> .	1287



1288	Alexander Scarlatos, Naiming Liu, Jaewook Lee,	Connor Shorten, Taghi M Khoshgoftaar, and Borko	1344
1289	Richard Baraniuk, and Andrew Lan. 2025. Train-	Furht. 2021. Text data augmentation for deep learn-	1345
1290	ing llm-based tutors to improve student learn-	ing. <i>Journal of big Data</i> , 8(1):101.	1346
1291	ing outcomes in dialogues. <i>arXiv preprint</i>		
1292	<i>arXiv:2503.06424</i> .		
1293	Torben Schmidt and Thomas Strasser. 2022. Artificial	Li Siyan, Teresa Shao, Zhou Yu, and Julia Hirschberg.	1347
1294	intelligence in foreign language learning and teach-	2024. <a href="#">EDEN: Empathetic dialogues for English</a>	1348
1295	ing: a call for intelligent practice. <i>Anglistik: Interna-</i>	<a href="#">learning</a> . In <i>Findings of the Association for Computa-</i>	1349
1296	<i>tional Journal of English Studies</i> , 33(1):165–184.	<i>tional Linguistics: EMNLP 2024</i> , pages 3492–3511,	1350
		Miami, Florida, USA. Association for Computational	1351
		Linguistics.	1352
1297	Robin Schmucker, Meng Xia, Amos Azaria, and Tom	Jiayang Song, Yuheng Huang, Zhehua Zhou, and Lei	1353
1298	Mitchell. 2024. Ruffle&riley: Insights from design-	Ma. 2024a. Multilingual blending: Llm safety	1354
1299	ing and evaluating a large language model-based con-	alignment evaluation with language mixture. <i>arXiv</i>	1355
1300	versational tutoring system. In <i>International Confer-</i>	<i>preprint arXiv:2407.07342</i> .	1356
1301	<i>ence on Artificial Intelligence in Education</i> , pages		
1302	75–90. Springer.	SeungWoo Song, Junghun Yuk, ChangSu Choi,	1357
		HanGyeol Yoo, HyeonSeok Lim, KyungTae Lim,	1358
1303	Johannes Schneider, Bernd Schenk, and Christina	and Jungyeul Park. 2025. <a href="#">Unified automated essay</a>	1359
1304	Niklaus. 2023. Towards llm-based autograd-	<a href="#">scoring and grammatical error correction</a> . In <i>Find-</i>	1360
1305	ing for short textual answers. <i>arXiv preprint</i>	<i>ings of the Association for Computational Linguistics:</i>	1361
1306	<i>arXiv:2309.11508</i> .	<i>NAACL 2025</i> , pages 4412–4426, Albuquerque, New	1362
		Mexico. Association for Computational Linguistics.	1363
1307	Kathrin Seßler, Maurice Fürstenberg, Babette Bühler,	Yixiao Song, Kalpesh Krishna, Rajesh Bhatt, Kevin	1364
1308	and Enkelejda Kasneci. 2024. Can ai grade your	Gimpel, and Mohit Iyyer. 2024b. <a href="#">GEE! grammar</a>	1365
1309	essays? a comparative analysis of large language	<a href="#">error explanation with large language models</a> . In	1366
1310	models and teacher ratings in multidimensional essay	<i>Findings of the Association for Computational Lin-</i>	1367
1311	scoring. <i>arXiv preprint arXiv:2411.16337</i> .	<i>guistics: NAACL 2024</i> , pages 754–781, Mexico City,	1368
		Mexico. Association for Computational Linguistics.	1369
1312	Zekai Shao, Siyu Yuan, Lin Gao, Yixuan He, Deqing	Shashank Sonkar, Kangqi Ni, Sapana Chaudhary, and	1370
1313	Yang, and Siming Chen. 2025. Unlocking scientific	Richard Baraniuk. 2024. <a href="#">Pedagogical alignment of</a>	1371
1314	concepts: How effective are llm-generated analogies	<a href="#">large language models</a> . In <i>Findings of the Associa-</i>	1372
1315	for student understanding and classroom practice?	<i>tion for Computational Linguistics: EMNLP 2024</i> ,	1373
1316	<i>arXiv preprint arXiv:2502.16895</i> .	pages 13641–13650, Miami, Florida, USA. Associa-	1374
		tion for Computational Linguistics.	1375
1317	Shikhar Sharma, Manas Mhasakar, Apurv Mehra,	Sinan Sonlu, Bennie Bendiksen, Funda Durupinar,	1376
1318	Utkarsh Venaik, Ujjwal Singhal, Dhruv Kumar,	and Uğur Güdükbay. 2024. The effects of em-	1377
1319	and Kashish Mittal. 2024. Comuniqa: Exploring	bodiment and personality expression on learning	1378
1320	large language models for improving english speak-	in llm-based educational agents. <i>arXiv preprint</i>	1379
1321	ing skills. In <i>Proceedings of the 7th ACM SIG-</i>	<i>arXiv:2407.10993</i> .	1380
1322	<i>CAS/SIGCHI Conference on Computing and Sustain-</i>		
1323	<i>able Societies</i> , pages 256–267.	Maja Stahl, Leon Biermann, Andreas Nehring, and Hen-	1381
1324	Lingfeng Shen, Weiting Tan, Sihao Chen, Yunmo Chen,	ning Wachsmuth. 2024a. <a href="#">Exploring LLM prompting</a>	1382
1325	Jingyu Zhang, Haoran Xu, Boyuan Zheng, Philipp	<a href="#">strategies for joint essay scoring and feedback gen-</a>	1383
1326	Koehn, and Daniel Khashabi. 2024a. The language	<a href="#">eration</a> . In <i>Proceedings of the 19th Workshop on</i>	1384
1327	barrier: Dissecting safety challenges of llms in multi-	<i>Innovative Use of NLP for Building Educational Ap-</i>	1385
1328	lingual contexts. <i>arXiv preprint arXiv:2401.13136</i> .	<i>plications (BEA 2024)</i> , pages 283–298, Mexico City,	1386
		Mexico. Association for Computational Linguistics.	1387
1329	Shuanghong Shen, Qi Liu, Zhenya Huang, Yonghe	Maja Stahl, Leon Biermann, Andreas Nehring, and Hen-	1388
1330	Zheng, Minghao Yin, Minjuan Wang, and Enhong	ning Wachsmuth. 2024b. Exploring llm prompting	1389
1331	Chen. 2024b. A survey of knowledge tracing: Mod-	strategies for joint essay scoring and feedback gener-	1390
1332	els, variants, and applications. <i>IEEE Transactions on</i>	ation. <i>arXiv preprint arXiv:2404.15845</i> .	1391
1333	<i>Learning Technologies</i> .		
1334	Yao Shi, Rongkeng Liang, and Yong Xu. 2025. Educa-	John Stamper, Ruiwei Xiao, and Xinying Hou. 2024.	1392
1335	tionq: Evaluating llms’ teaching capabilities through	Enhancing llm-based feedback: Insights from intelli-	1393
1336	multi-agent dialogue framework. <i>arXiv preprint</i>	gent tutoring systems and the learning sciences. In	1394
1337	<i>arXiv:2504.14928</i> .	<i>International Conference on Artificial Intelligence in</i>	1395
		<i>Education</i> , pages 32–43. Springer.	1396
1338	Mostafa Faghih Shojaei, Rahul Gulati, Benjamin A	Danny D Steinberg and Natalia V Sciarini. 2013. <i>An</i>	1397
1339	Jasperson, Shangshang Wang, Simone Cimolato,	<i>introduction to psycholinguistics</i> . Routledge.	1398
1340	Dangli Cao, Willie Neiswanger, and Krishna		
1341	Garikipati. 2025. Ai-university: An llm-based plat-		
1342	form for instructional alignment to scientific class-		
1343	rooms. <i>arXiv preprint arXiv:2504.08846</i> .		



1507	and beyond for grammatical error correction. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 10161–10175, Singapore. Association for Computational Linguistics.	
1508		
1509		
1510		
1511	Jingheng Ye, Shang Qin, Yinghui Li, Xuxin Cheng, Libo Qin, Hai-Tao Zheng, Peng Xing, Zishan Xu, Guo Cheng, and Zhao Wei. 2024. Exgcec: A benchmark of edit-wise explainable chinese grammatical error correction. <i>arXiv preprint arXiv:2407.00924</i> .	
1512		
1513		
1514		
1515		
1516	Yong Yu, Xiaosheng Si, Changhua Hu, and Jianxun Zhang. 2019. A review of recurrent neural networks: Lstm cells and network architectures. <i>Neural computation</i> , 31(7):1235–1270.	
1517		
1518		
1519		
1520	Murong Yue, Wijdane Mifdal, Yixuan Zhang, Jennifer Suh, and Ziyu Yao. 2024. Mathvc: An llm-simulated multi-character virtual classroom for mathematics education. <i>arXiv preprint arXiv:2404.06711</i> .	
1521		
1522		
1523		
1524	Weihaio Zeng, Lulu Zhao, Keqing He, Ruotong Geng, Jingang Wang, Wei Wu, and Weiran Xu. 2023. Seen to unseen: Exploring compositional generalization of multi-attribute controllable dialogue generation. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 14179–14196, Toronto, Canada. Association for Computational Linguistics.	
1525		
1526		
1527		
1528		
1529		
1530		
1531		
1532	Zhiyuan Zeng, Jiatong Yu, Tianyu Gao, Yu Meng, Tanya Goyal, and Danqi Chen. 2024. Evaluating large language models at evaluating instruction following. In <i>The Twelfth International Conference on Learning Representations</i> .	
1533		
1534		
1535		
1536		
1537	Daochen Zha, Zaid Pervaiz Bhat, Kwei-Herng Lai, Fan Yang, Zhimeng Jiang, Shaochen Zhong, and Xia Hu. 2023. Data-centric artificial intelligence: A survey. <i>ACM Computing Surveys</i> .	
1538		
1539		
1540		
1541	Bojun Zhan, Teng Guo, Xueyi Li, Mingliang Hou, Qianru Liang, Boyu Gao, Weiqi Luo, and Zitao Liu. 2024. Knowledge tracing as language processing: A large-scale autoregressive paradigm. In <i>International Conference on Artificial Intelligence in Education</i> , pages 177–191. Springer.	
1542		
1543		
1544		
1545		
1546		
1547	Hanqing Zhang, Haolin Song, Shaoyu Li, Ming Zhou, and Dawei Song. 2023. A survey of controllable text generation using transformer-based pre-trained language models. <i>ACM Computing Surveys</i> , 56(3):1–37.	
1548		
1549		
1550		
1551		
1552	Ningyu Zhang, Yunzhi Yao, Bozhong Tian, Peng Wang, Shumin Deng, Mengru Wang, Zekun Xi, Shengyu Mao, Jintian Zhang, Yuansheng Ni, and 1 others. 2024a. A comprehensive study of knowledge editing for large language models. <i>arXiv preprint arXiv:2401.01286</i> .	
1553		
1554		
1555		
1556		
1557		
1558	Ruqing Zhang, Jiafeng Guo, Lu Chen, Yixing Fan, and Xueqi Cheng. 2021. A review on question generation from natural language text. <i>ACM Transactions on Information Systems (TOIS)</i> , 40(1):1–43.	
1559		
1560		
1561		
	Zheyuan Zhang, Daniel Zhang-Li, Jifan Yu, Linlu Gong, Jinchang Zhou, Zhiyuan Liu, Lei Hou, and Juanzi Li. 2024b. Simulating classroom education with llm-empowered agents. <i>arXiv preprint arXiv:2406.19226</i> .	1562
		1563
		1564
		1565
		1566
	Daniel Zhang-Li, Zheyuan Zhang, Jifan Yu, Joy Lim Jia Yin, Shangqing Tu, Linlu Gong, Haohua Wang, Zhiyuan Liu, Huiqin Liu, Lei Hou, and 1 others. 2024. Awakening the slides: A tuning-free and knowledge-regulated ai tutoring system via language model co-ordination. <i>arXiv preprint arXiv:2409.07372</i> .	1567
		1568
		1569
		1570
		1571
		1572
	Ying Zheng, Xueyi Li, Yaying Huang, Qianru Liang, Teng Guo, Mingliang Hou, Boyu Gao, Mi Tian, Zitao Liu, and Weiqi Luo. 2024. Automatic lesson plan generation via large language models with self-critique prompting. In <i>International Conference on Artificial Intelligence in Education</i> , pages 163–178. Springer.	1573
		1574
		1575
		1576
		1577
		1578
		1579
	Wanjun Zhong, Lianghong Guo, Qiqi Gao, He Ye, and Yanlin Wang. 2024. Memorybank: Enhancing large language models with long-term memory. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pages 19724–19731.	1580
		1581
		1582
		1583
		1584
	Ruiyang Zhou, Lu Chen, and Kai Yu. 2024. Is LLM a reliable reviewer? a comprehensive evaluation of LLM on automatic paper reviewing tasks. In <i>Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)</i> , pages 9340–9351, Torino, Italia. ELRA and ICCL.	1585
		1586
		1587
		1588
		1589
		1590
		1591
	Yizhou Zhou, Mengqiao Zhang, Yuan-Hao Jiang, Xinyu Gao, Naijie Liu, and Bo Jiang. 2025. A study on educational data analysis and personalized feedback report generation based on tags and chatgpt. <i>arXiv preprint arXiv:2501.06819</i> .	1592
		1593
		1594
		1595
		1596
	Xinlin Zhuang, Hongyi Wu, Xinshu Shen, Peimin Yu, Gaowei Yi, Xinhao Chen, Tu Hu, Yang Chen, Yuyei Ren, Yadong Zhang, Youqi Song, Binxuan Liu, and Man Lan. 2024. TOREE: Evaluating topic relevance of student essays for Chinese primary and middle school education. In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 5749–5765, Bangkok, Thailand. Association for Computational Linguistics.	1597
		1598
		1599
		1600
		1601
		1602
		1603
		1604
		1605



## A Literature Review

We provide an overview of LLM-centric research of English Education presented in Figure 4.

## B Four Phases of Research Roadmap

### Stage 1: Rule-based Models (1960s–1990s).

Early solutions relied on handcrafted linguistic rules to process language in tightly constrained scenarios (Grosan et al., 2011; C Angelides and Garcia, 1993). Classical platforms like PLATO (Hart, 1981) and Systran (Toma, 1977) operated effectively for highly structured tasks (e.g., grammar drills) but struggled with complex, context-dependent interactions.

### Stage 2: Statistical Models (1990s–2010s).

With the increased availability of digitized corpora, methods such as the early version of Google Translate (Och, 2006) and Dragon NaturallySpeaking (Blair, 1997) pioneered statistical pattern mining. These approaches leveraged large datasets to infer linguistic rules probabilistically, improving scalability yet still lacking deeper semantic understanding.

**Stage 3: Neural Models (2010s–2020s).** The advent of deep learning architectures (e.g., RNNs (Yu et al., 2019) and Transformers (Vaswani, 2017)) enabled more robust context modeling, sparking transformative applications like Grammarly (Fitria, 2021) and Duolingo (Vesselinov and Grego, 2012). These systems offered enhanced personalization and feedback, significantly augmenting learners’ writing and reading comprehension.

**Stage 4: Large Language Models (2020s–Present).** Nowadays, various LLMs (e.g., ChatGPT (Achiam et al., 2023)) combine massive pre-training with instruction tuning, achieving impressive results in multi-turn dialogue, individualized scaffolding, and multimodal integration. Tools such as Khanmigo (Anand, 2023) demonstrate LLMs’ potential for real-time conversational practice, dynamic content creation, and inclusive educational support at scale.

## C Challenges and Future Directions

While we posit that LLMs have the potential to revolutionize English Education, realizing their full promise requires addressing key challenges. This section offers a concise overview of these chal-

lenges, followed by directions that could guide future research and deployment.

### Ensuring Reliability and Mitigating Hallucina-

tions. LLMs may produce hallucinations (Huang et al., 2023) that can mislead learners and undermine pedagogical goals. This risk intensifies in high-stakes educational environments, where trust and correctness are paramount. Future directions include enhancing data quality and diversity for training (Long et al., 2024), developing techniques to integrate LLM outputs with structured domain knowledge and pedagogical rules, and employing rigorous automated and human-in-the-loop validation mechanisms to minimize such detrimental outcomes and improve the factual grounding of LLM-generated educational content.

### Addressing Bias and Ethical Considerations.

As LLMs inherit biases from their training data, these systems may produce culturally insensitive, stereotypical, or unfair responses, potentially harming students from diverse linguistic and sociocultural backgrounds. Moreover, significant privacy concerns emerge when collecting and using learner data to personalize instruction, particularly for K-12 students. Future research must focus on developing robust governance frameworks, transparent documentation of data sources and model behaviors, and advanced bias detection and mitigation strategies (Borah and Mihalcea, 2024; He and Li, 2024) to ensure that LLM-based tools for English Education are equitable, fair, and uphold stringent data protection standards.

### Aligning With Pedagogical Principles.

LLMs excel at generating fluent language but often lack deep pedagogical alignment, particularly for tasks requiring developmental sensitivity, learner motivation strategies, or differentiated instruction tailored to individual learning needs. Their general-purpose nature means they do not inherently account for established language acquisition theories or specific curricular standards (Razafinirina et al., 2024). A crucial future direction is the development of methodologies to better imbue LLMs with pedagogical intelligence. This includes co-designing LLM applications with educators, fine-tuning models on high-quality pedagogical interaction data, and creating architectures that can dynamically adapt to learners’ cognitive states and developmental needs in English language learning.

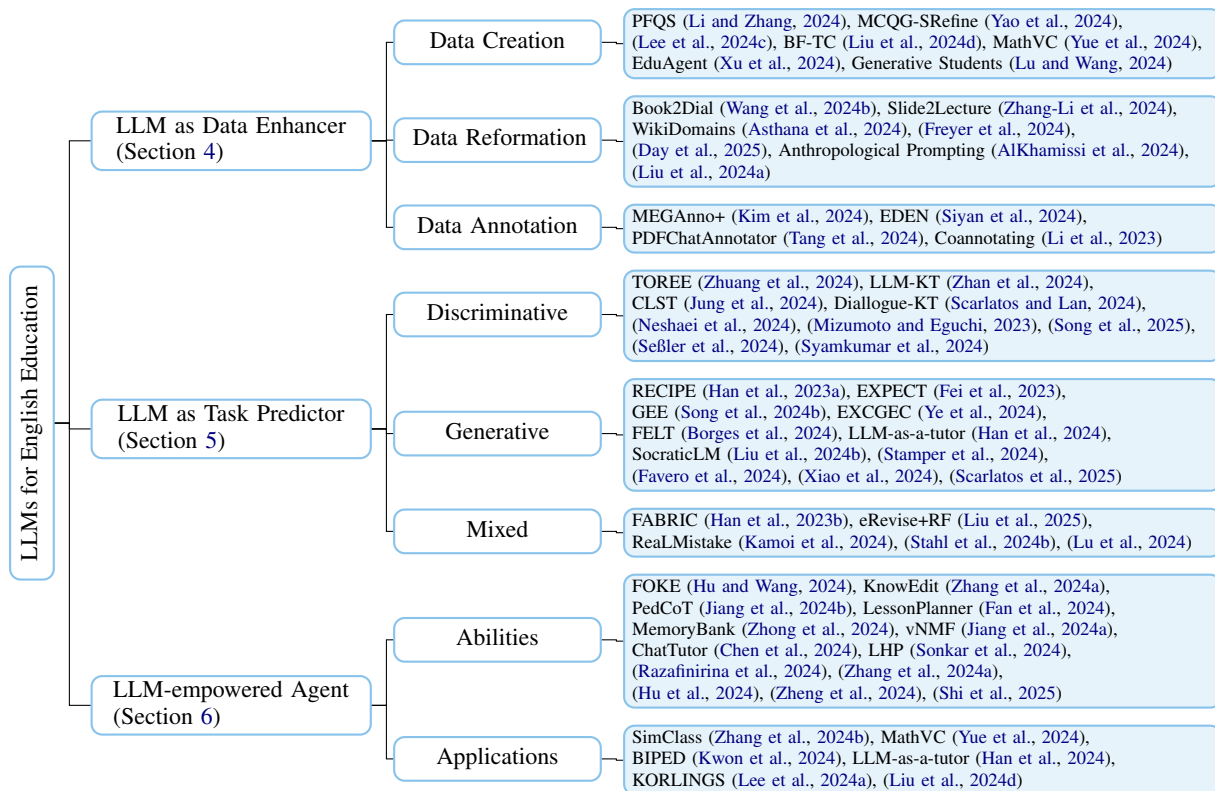


Figure 4: An overview of LLM-centric research of FLE.

## Establishing Robust Evaluation Frameworks.

A significant challenge in leveraging LLMs for English Education is the current lack of widely accepted and easily implementable evaluation frameworks to assess the quality of LLM-based teaching interactions and outcomes. Existing metrics often focus on linguistic correctness or task completion (Tan et al., 2024; Macina et al., 2025) rather than pedagogical efficacy or impact on learning (Chiang et al., 2024). Future work should prioritize the development of standardized evaluation methodologies, including comprehensive benchmarks and nuanced metrics that capture both the accuracy of linguistic information and the pedagogical value of LLM interventions. This will be essential for comparing different systems and guiding iterative improvements.

## Integrating with Standardized Educational Frameworks.

English language learning is often governed by established standards and frameworks, such as the Common European Framework of Reference for Languages (CEFR) or Common Core State Standards (CCSS). For LLM-based tools to be truly effective and gain acceptance, their outputs and interaction patterns should align with these existing frameworks. Future technical de-

velopment should focus on enabling LLMs to reference, interpret, and operate consistently within these standards (Nicholls et al., 2024; Imperial et al., 2024). This includes generating proficiency-level-appropriate content, providing feedback that corresponds to specific framework descriptors, and assisting learners in achieving standardized learning objectives, thereby enhancing usability, conformity, and trustworthiness among educators and learners.

## Fostering Human-AI Collaboration in Pedagogy.

While LLMs offer transformative potential, it is unlikely they will completely replace human teachers in English Education in the foreseeable future. Instead, the most promising path involves developing sophisticated human-AI collaborative educational technologies (KIM et al., 2025). Future research should explore how LLMs can best function as assistive tools that augment, rather than supplant, the capabilities of human educators (Shojaei et al., 2025). This includes designing intuitive interfaces for teachers to guide, customize, and oversee LLM-driven activities, investigating teachers' perspectives on integrating LLMs into their practice, and defining technical benchmarks for when an LLM possesses sufficient acquired skills to reliably as-

sist teachers. The focus must be on a synergistic model where LLMs handle scalable tasks while human teachers provide the crucial elements of empathy, nuanced understanding, and holistic student development.

systems (Wu et al., 2022), can help minimize inaccuracies (Ho et al., 2024). Additionally, the fine-tuning capabilities of LLMs ensure adaptability, supporting diverse and inclusive learning experiences (Lee et al., 2024b).

## D Alternative Views

While this paper supports using LLMs in English Education, it is essential to consider alternative perspectives. Below, we discuss two key opposing views and provide counterarguments.

### D.1 Task-Specific or Language-Specific Models as Better Alternatives

Some argue that specialized or language-specific models, including classical ML systems with carefully engineered features, can outperform general-purpose LLMs in narrowly defined tasks (e.g., phonetics or grammar drills (Fang et al., 2023)). By focusing on limited objectives, such models avoid the computational overhead and potential inaccuracies of LLMs, which aim to handle a broader range of inputs and contexts (Shen et al., 2024a).

**Counterargument.** While specialized models may excel in isolated tasks, they lack the flexibility required for comprehensive English Education, which involves cultural nuances, conversations, and evolving learner needs. In contrast, LLMs can be fine-tuned for specific goals while still offering broader linguistic competence (Song et al., 2024a). Additionally, relying on multiple specialized models can be resource-intensive, whereas a well-configured LLM provides a unified framework that balances specialization and scalability.

### D.2 Concerns About Over-Reliance on LLMs

Critics warn that over-reliance on LLMs may lead to problems such as generating misleading outputs (Nahar et al., 2024), reducing human interaction, and over-standardizing teaching methods. These issues could undermine the interpersonal and motivational aspects of language learning.

**Counterargument.** These risks highlight the need for balanced integration rather than the replacement of human tutors. LLMs can complement educators by automating repetitive tasks, allowing teachers to focus on individualized support and motivation. Advances in AI safety, such as feedback loops (Tong et al., 2024) and human-in-the-loop