

# From Passive Tool to Socio-cognitive Teammate: A Conceptual Framework for Agentic AI in Human-AI Collaborative Learning

Lixiang Yan<sup>a,b,\*</sup>

<sup>a</sup>School of Education, Tsinghua University, Beijing, China

<sup>b</sup>Centre for Learning Analytics, Monash University, Clayton, Australia

## ARTICLE INFO

### Keywords:

Agentic AI  
Human-AI Collaboration  
Collaborative Learning  
Computer-Supported Collaborative Learning (CSCL)  
AI in Education

## ABSTRACT

The role of Artificial Intelligence (AI) in education is undergoing a rapid transformation, moving beyond its historical function as an instructional tool towards a new potential as an active participant in the learning process. This shift is driven by the emergence of agentic AI, autonomous systems capable of proactive, goal-directed action. However, the field lacks a robust conceptual framework to understand, design, and evaluate this new paradigm of human-AI interaction in learning. This paper addresses this gap by proposing a novel conceptual framework (the APCP framework) that charts the transition from AI as a tool to AI as a collaborative partner. We present a four-level model of escalating AI agency within human-AI collaborative learning: (1) the AI as an Adaptive Instrument, (2) the AI as a Proactive Assistant, (3) the AI as a Co-Learner, and (4) the AI as a Peer Collaborator. Grounded in sociocultural theories of learning and Computer-Supported Collaborative Learning (CSCL), this framework provides a structured vocabulary for analysing the shifting roles and responsibilities between human and AI agents. The paper further engages in a critical discussion of the philosophical underpinnings of collaboration, examining whether an AI, lacking genuine consciousness or shared intentionality, can be considered a true collaborator. We conclude that while AI may not achieve authentic phenomenological partnership, it can be designed as a highly effective functional collaborator. This distinction has significant implications for pedagogy, instructional design, and the future research agenda for AI in education, urging a shift in focus towards creating learning environments that harness the complementary strengths of both human and AI.

## 1. The Next Frontier of Educational AI

For decades, the integration of Artificial Intelligence in Education (AIED) has been a subject of intense research and development, promising to transform teaching and learning (Yan, Greiff, Teuber and Gašević, 2024; Giannakos, Azevedo, Brusilovsky, Cukurova, Dimitriadis, Hernandez-Leo, Järvelä, Mavrikis and Rienties, 2025; Chen, Zou, Xie, Cheng and Liu, 2022). Historically, this promise has been pursued primarily through the lens of individualization and efficiency. The predominant applications of AIED have been Intelligent Tutoring Systems (ITS) and adaptive learning platforms, which leverage AI to provide personalized instruction, real-time feedback, and customized learning pathways (Ouyang and Jiao, 2021; Kulik and Fletcher, 2016). These systems, often built on cognitive and mastery-learning principles, have demonstrated effectiveness in specific, well-defined domains (Kulik and Fletcher, 2016). However, they have also faced criticism for frequently replicating traditional, teacher-centric pedagogical models, where knowledge is transmitted to a passive learner (Ouyang and Jiao, 2021; Kulik and Fletcher, 2016). Consequently, the potential of AI to support more constructivist and socially-oriented modes of learning, such as collaborative learning, has remained largely unrealized (Zhou and Schofield, 2024).

The field of AIED and educational technology is now at a critical inflection point with the maturity of *agentic AI*. This new class of AI represents a fundamental paradigm shift. Unlike traditional AI systems that are largely reactive, agentic AI is defined by its autonomy, proactivity, and goal-driven behaviour, it can perceive its environment, reason about its goals, and execute complex, multi-step actions with limited human supervision (Sapkota, Roumeliotis and Karkee, 2025; Kamalov, Calonge, Smail, Azizov, Thadani, Kwong and Atif, 2025). This transition from a reactive tool to a proactive actor challenges the established roles and power dynamics within the learning process. The discourse must evolve from a focus on "using AI for learning" to one of "learning with AI," where the AI is not merely a resource but an active participant in the co-construction of knowledge.

ORCID(s): 0000-0003-3818-045X (L. Yan)

This emergent capability introduces a conceptual challenge. The existing models and frameworks used to understand AIED, which are largely built around the AI-as-tutor or AI-as-tool metaphor, are insufficient for capturing the nuances of a human-AI collaborative partnership (Chen et al., 2022; Ouyang and Jiao, 2021). The interaction is no longer simply about a user operating a piece of software; it becomes a dynamic interplay between two agents, one human, one artificial, that must coordinate their actions to achieve a shared objective (Yusuf, Money and Daylamani-Zad, 2025; Sapkota et al., 2025). This necessitates a move away from frameworks rooted in traditional Human-Computer Interaction (HCI), which prioritize usability and task performance, towards frameworks that can account for the complexities of collaboration, negotiation, and shared goals, drawing inspiration from the rich traditions of Computer-Supported Collaborative Learning (CSCL) and social psychology.

Foundational perspectives from scholars like Ben Shneiderman 2020 and Mutlu Cukurova 2025 offer essential guideposts for navigating this new terrain. Shneiderman's Human-Centered AI (HCAI) framework establishes a crucial design philosophy, advocating for systems that strategically combine high levels of human control and high levels of computer automation. The goal is to develop technologies that are demonstrably reliable, safe, and trustworthy, thereby enhancing human performance and creativity (Shneiderman, 2020). Complementing this, Cukurova's 2025 AIED-HCD framework offers a valuable educational typology, classifying AI's impact on human competence as either externalizing, internalizing, or extending cognition (Cukurova, 2025). Our proposed framework seeks to build directly upon these influential ideas, offering a more focused lens to operationalize their high-level visions within the specific, dynamic context of collaborative learning. While these models provide the "why" and the "what," our framework provides a "how" for designing the nuanced, moment-to-moment interactions with an agentic AI partner. A dedicated framework for AI agency is therefore needed to articulate the graduated roles an AI partner can inhabit, guiding the design of truly synergistic and trustworthy human-AI teams.

To harness the pedagogical potential of agentic AI, the field requires a new conceptual language to describe, design, and evaluate these nascent partnerships. This paper seeks to provide such a language. We propose a conceptual framework, the **APCP** (Adaptive instrument, Proactive assistant, Co-learner, Peer collaborator) framework, that moves beyond the simplistic tool-partner dichotomy to outline four distinct levels of AI agency in the context of human-AI collaborative learning. This framework offers a vocabulary for researchers, designers, and educators to articulate and navigate the evolving relationship between human learners and their increasingly capable artificial counterparts.

## 2. The Sociocultural Foundations of Collaborative Learning

To conceptualize the role of an agentic AI as a collaborator, it is first essential to establish a robust theoretical understanding of what collaboration entails. The field of learning sciences provides a deep and nuanced definition, rooted in sociocultural theory, that positions collaboration not merely as group work, but as a fundamental mechanism of human learning and cognitive development (Dillenbourg, 1999b).

Collaborative learning is broadly defined as an educational approach where two or more individuals learn together by working on a joint task to achieve a common goal (Dillenbourg, 1999b). This perspective is heavily influenced by the work of Vygotsky 1978, who posited that learning is an inherently social activity, mediated by interaction with the social environment. Knowledge is not seen as an objective entity to be transferred from an expert to a novice, but as something that is actively co-constructed by learners through dialogue, negotiation, and the shared use of tools (Bruffee, 1999; Roschelle and Teasley, 1995). This social constructivist view frames learning as a process of *intersubjective meaning making* (Dillenbourg, 1999a), where participants strive to build and maintain a shared understanding of the problem and their joint activity (Dillenbourg, 1999b; Roschelle and Teasley, 1995). As Bruffee 1999 describes it, collaborative learning "creates conditions in which students can negotiate the boundaries between the knowledge communities they belong to and the one that the professor belongs to."

A central concept in this tradition is Vygotsky's Zone of Proximal Development (ZPD), which refers to the gap between what a learner can achieve independently and what they can accomplish with guidance from a more capable peer or instructor (Vygotsky, 1978). The ZPD underscores the importance of the collaborator as a scaffold, enabling learners to extend their capabilities beyond what they could achieve alone. This collaborative interaction is not merely about receiving assistance; it serves as a driving force for cognitive development.

From these theoretical foundations, several core principles of effective collaboration can be derived. It requires more than just co-presence; it demands the mutual engagement of participants in a coordinated effort to solve the problem together (Bruffee, 1999; Roschelle and Teasley, 1995). Participants must actively facilitate and encourage one another's contributions in a process of promotive interaction (Johnson and Johnson, 1991). The cornerstone of this

process is the establishment of a shared goal and intersubjectivity, which involves building a joint problem space and a shared understanding that allows participants to coordinate their actions and build common knowledge (Dillenbourg, 1999b; Roschelle and Teasley, 1995; Dillenbourg, 1999a). Finally, while the effort is collective, each member remains responsible for their individual accountability, ensuring that all participants are actively engaged rather than passively observing (Johnson and Johnson, 1991).

It is useful to distinguish this rich, philosophical view of *collaborative* learning from the more structured, method-oriented approach of *cooperative* learning. While both involve small group work, cooperative learning tends to be more highly structured by the instructor, with a greater emphasis on division of labour and individual accountability for specific outcomes (Johnson and Johnson, 1991). Collaborative learning, in contrast, is often more loosely structured, emphasizing the process of knowledge negotiation and the social construction of meaning itself (Bruffee, 1999). For the purposes of this paper, which seeks to explore the potential for a deep, partnered relationship with AI, the more demanding, philosophical definition of collaboration serves as the ideal benchmark.

The field of CSCL emerged to study how technology can mediate, facilitate, and scaffold these complex collaborative processes (Dillenbourg, 1999b; Roschelle and Teasley, 1995; Lehtinen, 2003). CSCL research demonstrates that technology is not a neutral conduit; its design can either support or hinder effective collaboration (Jeong and Hmelo-Silver, 2016; Lehtinen, 2003). The success of a CSCL environment depends critically on the pedagogical design that structures the interaction, not on the technological features alone (Jeong and Hmelo-Silver, 2016; Lehtinen, 2003).

This theoretical grounding reveals a crucial point for the present discussion. The very definition of collaboration is laden with deeply human-centric concepts: "shared understanding," "negotiation of meaning," "intersubjectivity," and "mutual engagement." These concepts presuppose the existence of conscious subjects who possess beliefs, intentions, and perspectives that can be shared and negotiated. This presents a formidable conceptual challenge for any non-human agent. A purely technical or behavioural replication of collaborative actions may fail to capture the pedagogical essence of the process if it cannot engage with these underlying social and cognitive dynamics. This inherent tension between the human-centric nature of collaboration and the artificial nature of AI necessitates the critical analysis in Section 5, forcing us to question whether we are aiming to build an AI that is a *true* collaborator or one that is a *functionally effective* one.

### 3. Agentic AI: From Reactive Tool to Proactive Partner

To argue for a new framework, it is necessary to establish that agentic AI represents a genuinely new category of educational technology, distinct from its predecessors. Its unique characteristics are what create the need for a new conceptual model of interaction. Agentic AI is defined as an autonomous system that can perceive its environment, reason about complex goals, and act independently to achieve them with minimal human supervision (Park, O'Brien, Cai, Morris, Liang and Bernstein, 2023; Durante, Huang, Wake, Gong, Park, Sarkar, Taori, Noda, Terzopoulos, Choi et al., 2024; Dai, Ke, Pan, Moon and Liu, 2024). The term "agentic" signifies this capacity for agency, the ability to make independent, purposeful, and context-aware decisions rather than simply following pre-defined rules or direct commands (Park et al., 2023; Sapkota et al., 2025; Kamalov et al., 2025).

The operation of an agentic system can be understood through a continuous perception-reasoning-action loop (Wang, Ma, Feng, Zhang, Yang, Zhang, Chen, Tang, Chen, Lin et al., 2024; Park et al., 2023; Durante et al., 2024). The process begins with *Perception*, where the agent collects real-time data from its environment, such as user interactions, database content, or the state of other software systems (Wang et al., 2024; Durante et al., 2024). Following this, the agent enters a phase of *Reasoning and Goal Setting*, using capabilities like natural language processing (NLP) to interpret user intent and formulate a strategy, often by breaking down high-level goals into a sequence of smaller, actionable sub-tasks (Wang et al., 2024; Park et al., 2023). This leads to *Decision-Making and Execution*, where the agent evaluates potential actions, selects the optimal one, and executes it by interacting with external tools, responding to the user, or orchestrating other agents (Wang et al., 2024; Park et al., 2023). Crucially, the loop is completed through *Learning and Adaptation*, as the agent gathers feedback from the outcomes of its actions, using techniques like reinforcement learning to evaluate its performance and refine its strategies over time (Wang et al., 2024; Park et al., 2023; Durante et al., 2024).

A critical distinction must be made between agentic AI and the more familiar generative AI (e.g., ChatGPT). While agentic systems often leverage generative models for their reasoning and communication capabilities, their function is fundamentally different. Generative AI is built to create content, it produces text, images, or code in response to a user's prompt (Xi, Chen, Guo, He, Ding, Hong, Zhang, Wang, Jin, Zhou et al., 2025). Its role is primarily reactive. Agentic

**Table 1**  
The Evolution of AI's Role in Education

Dimension	Intelligent Tutoring System (ITS)	Generative AI (e.g., Chat-GPT)	Agentic AI
<b>Primary Function</b>	Provide adaptive instruction and feedback within a defined domain	Generate novel content (text, code, images) based on user prompts	Autonomously perform multi-step tasks to achieve educational goals
<b>Locus of Control</b>	System-driven within a pre-defined curriculum; learner follows a path	Human-driven via prompts; AI is reactive	Shared or AI-driven; AI exhibits autonomy and proactivity
<b>Interaction Model</b>	Structured, question-answer, feedback loops	Conversational, prompt-response	Goal-oriented, can initiate actions, orchestrate tools and other agents
<b>Adaptability</b>	Adapts content difficulty and sequencing based on learner performance	Adapts responses based on conversational context	Adapts its entire strategy and action plan based on environmental feedback and outcomes
<b>Core Learning Theory</b>	Primarily cognitive; mastery learning	None intrinsic (can support multiple pedagogical approaches)	Potential for social constructivist and collaborative learning paradigms

AI, in contrast, is built to act. It uses the outputs of generative models as part of a larger process to autonomously plan and execute tasks to achieve a goal (Sapkota et al., 2025; Kamalov et al., 2025). For example, a generative AI can draft a set of quiz questions based on a chapter of a textbook; an agentic AI can draft the questions, analyze students' past performance data to tailor difficulty levels, automatically upload the quiz to the learning management system, notify the relevant student group, and schedule adaptive follow-up activities for learners who struggle with specific concepts.

This proactive, goal-oriented nature also distinguishes agentic AI from traditional ITS. An ITS provides personalized feedback and adapts the difficulty of content, but it operates within a highly structured, pre-defined curriculum and pedagogical model (Ouyang and Jiao, 2021; Kulik and Fletcher, 2016). It is domain-specific and reactive to student inputs. An agentic AI, conversely, can operate in unstructured, open-ended environments, set its own sub-goals, and dynamically alter its entire strategy based on the flow of the interaction (Sapkota et al., 2025; Kamalov et al., 2025). The very limitations often cited in ITS research, such as a failure to support constructivist, inquiry-based, or collaborative learning (Kulik and Fletcher, 2016), are precisely the areas where agentic AI holds the most transformative potential. Table 1 clarifies these distinctions, illustrating the paradigm shift that agentic AI represents.

## 4. A Framework for Agentic AI in Human-AI Collaborative Learning

The emergence of agentic AI necessitates a structured way to conceptualize its role in learning. Drawing on existing models of human-AI interaction (Shneiderman, 2020; Cukurova, 2025) and taxonomies of AI autonomy (Bradshaw, Hoffman, Woods and Johnson, 2013; Endsley, 2017), but tailoring them specifically to the pedagogical context of collaborative learning, this paper proposes a four-level framework: the APCP framework (Figure 1). This framework describes a continuum of escalating AI agency, where each level represents a distinct configuration of roles, responsibilities, and interaction dynamics between the human learner and the AI agent. It provides a vocabulary for describing, designing, and evaluating these new learning partnerships.

### 4.1. Level 1: The AI as an Adaptive Instrument

At this foundational level, the AI functions as a highly sophisticated, yet fundamentally passive, instrument. All significant cognitive and intentional agency resides exclusively with the human learner. The AI does not initiate tasks, make independent decisions, or pursue its own goals. Its actions are direct, deterministic responses to explicit human

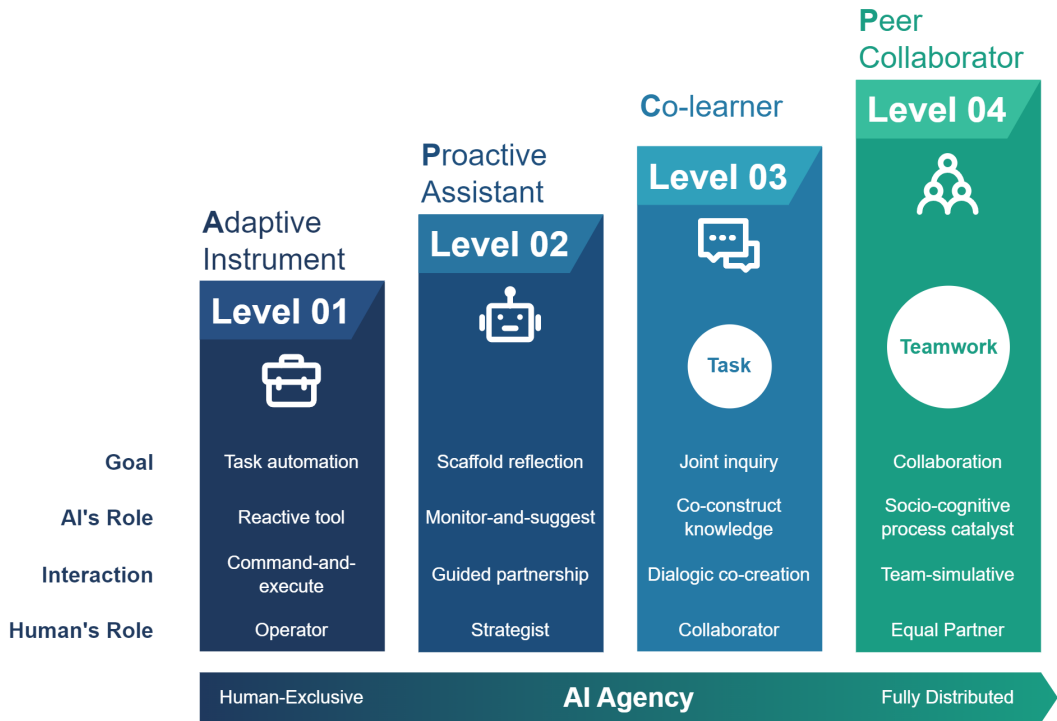


Figure 1: The APCP Framework: A Conceptual Framework for AI Agency in Human-AI Collaborative Learning

commands. This model aligns with the “user as an operator” or “copilot” metaphor, where the AI is an extension of the user’s will, available for on-demand support (Bradshaw et al., 2013). From a Vygotskian perspective, the AI acts as a powerful mediating artifact, a psychological tool that can augment and reshape cognition but does not participate in the cognitive process as a partner (Vygotsky, 1978). The human is the sole planner, strategist, and decision-maker, while the AI’s role is one of reactive execution.

The interaction at this level is characterized by a master-servant dynamic. The learner directs, and the AI performs. For instance, in a collaborative science project, two students might be analyzing a dataset on climate change. One student could issue a command: Generate a time-series plot of average global temperature from 1950 to 2020 and overlay a 10-year moving average. The AI executes this command precisely. Its “adaptivity” is limited to its capacity to parse the request and, perhaps, adjust the visualization’s complexity based on a rudimentary user model (e.g., simplifying the graph if the user’s prior queries were very basic). However, the impetus for action, the choice of visualization, and the interpretation of the output remain entirely human responsibilities.

The primary pedagogical value of Level 1 agentic AI is its potential to reduce extraneous cognitive load (Sweller, 2010). By automating laborious, low-level tasks such as data formatting, calculation, or information retrieval, it frees the learner’s cognitive resources to focus on higher-order thinking skills like analysis, synthesis, and evaluation. This reflects the traditional use of technology in CSCL environments, where the technology’s primary function is to facilitate human-to-human interaction and knowledge construction, rather than to act as a collaborator itself (O’Malley, 2012). Consequently, designing a Level 1 agentic AI requires a robust natural language processing (NLP) front-end to interpret user commands and a powerful back-end to execute a defined set of tasks. The focus is on reliability, speed, and accuracy of execution. The primary limitation of this level is that the AI provides no scaffolding for the collaborative process itself. It can support the task but not the teamwork. It cannot prompt for reflection, challenge a flawed assumption, or help resolve a disagreement between human collaborators.

Empirical evidence demonstrates that even at Level 1, where AI functions purely as a reactive tool, its targeted support can meaningfully enhance collaborative learning outcomes. In a quasi-experimental study in programming education, Wang, Wang, Chen, Liu, Bao and Xu (2025) found that an AI-Agent-supported collaborative learning model, in which the AI acted solely on explicit student prompts, significantly improved university students’



learning achievement, self-efficacy, and interest while lowering mental effort compared to a traditional CSCL control group. Similarly, Wei, Wang, Lee and Liu (2025) examined a 20-week digital storytelling project in which university teams used generative AI tools such as ChatGPT, Midjourney, and Runway for idea generation and content drafting. Although the AI offered no proactive scaffolding, students reported that its on-demand outputs enhanced collaborative problem-solving and team creativity, especially in generating novel ideas and improving user experience. These findings highlight that Level 1 agentic AI, while lacking initiative or strategic input, can still reduce cognitive load and amplify group productivity, thereby creating more space for human collaborators to focus on higher-order joint thinking.

#### 4.2. Level 2: The AI as a Proactive Assistant

The second level marks a significant shift: the AI gains a limited and bounded form of proactive agency. It transcends its role as a passive instrument to become an active assistant that can monitor the learning context and offer unsolicited support. This is a move towards partial autonomy, where the AI functions as a decision-support system (Endsley, 2017). It can augment the human learner's perception and attention by highlighting salient information or potential issues that might otherwise be overlooked (Kirsh, 2009). Crucially, however, the human learner retains ultimate strategic control and veto power. The human's role evolves from a direct operator to a strategist and reviewer, while the AI's role is to anticipate needs and provide opportune, actionable suggestions.

The interaction dynamic becomes a guided partnership. The AI interjects with suggestions, but the human must approve them, embodying what Shneiderman terms an "AI-assisted full human control mode" (Shneiderman, 2020). Consider a student drafting an argumentative essay on economic policy. The AI, monitoring the text in real-time, might detect a potential logical fallacy. Instead of waiting to be asked, it could proactively highlight the problematic sentence and display a message: "This statement appears to be a hasty generalization. The evidence you've presented so far only supports this claim for developed nations. Would you like me to find sources that discuss the policy's effect in emerging economies?"

The pedagogical purpose here is to scaffold metacognition and self-regulation. The AI's proactive prompts encourage the learner to reflect on their own work, identify weaknesses, and consider alternative perspectives. It acts as a cognitive nudge, pushing the learner toward more rigorous thinking without dictating the final product. For collaborative tasks, a Level 2 agentic AI could monitor the dialogue between two students and gently intervene (e.g., "It seems you both agree on the solution, but have you considered this potential counter-argument from the literature?"). This helps the group avoid premature consensus and encourages deeper critical inquiry. Therefore, building a Level 2 agentic AI is considerably more complex. It requires sophisticated user modeling to understand the learner's goals and current state of knowledge. It also needs a carefully designed "interruption protocol" to ensure its proactive suggestions are helpful rather than intrusive or disruptive to the learner's flow. The key design challenge is balancing proactivity with deference to human authority. A major limitation is that the AI's proactivity is typically constrained to pre-defined rules and triggers; it can identify known problems but cannot engage in novel, creative problem-solving with the learner. Its role remains one of support, not co-creation.

Proactive engagement by AI, when bounded and context-aware, is increasingly demonstrable and beneficial in collaborative and data-rich learning contexts. In a recent randomised controlled trial, Yan, Martinez-Maldonado, Jin, Echeverria, Milesi, Fan, Zhao, Alfredo, Li and Gašević (2025b) investigated the effects of generative AI agents with scaffolding on students' comprehension of complex visual learning analytics. The study compared three conditions, passive agents, proactive agents that employed scaffolding questions, and standalone scaffolding, across 117 higher education students. Results showed that proactive GenAI agents significantly improved student comprehension compared to both passive agents and standalone scaffolding, with these benefits persisting beyond the intervention. The proactive agents' ability to detect learner needs and pose timely, targeted questions aligns closely with Level 2 design principles, as they enhance metacognitive engagement without taking over the cognitive process. Similarly, Pu, Lazaro, Arawjo, Xia, Xiao, Grossman and Chen (2025) evaluated Codellaborator, a proactive AI programming assistant that monitored the coding context and initiated suggestions when appropriate. In a within-subject study (N = 18), Codellaborator improved programming efficiency compared to a prompt-only condition, though poorly timed interventions could disrupt workflow; interface variants with presence indicators and richer context mitigated this issue, illustrating the design challenge of balancing proactivity with learner control. Together, these findings illustrate that Level 2 agentic AI, by anticipating learner needs and intervening judiciously, can scaffold reflection, reduce oversight, and enhance collaborative task performance, while still preserving human authority and avoiding intrusive behavior.

### 4.3. Level 3: The AI as a Co-Learner

At this level, the relationship evolves into a more symmetrical, dialogic partnership. The AI is no longer merely a support system but a co-learner capable of tackling substantive parts of a problem in parallel with the human. The locus of agency is now shared and negotiated. This model aligns with the “user as a collaborator” paradigm, defined by rich, reciprocal communication and a shared cognitive workspace (Bradshaw et al., 2013). The human’s role shifts to that of a true collaborator and, at times, a mentor who might need to “teach” the AI to refine its approach. The AI, in turn, can model learning processes, articulate its own functional “uncertainty,” and engage in the co-construction of knowledge.

The interaction mirrors human peer collaboration. A central feature is the AI’s ability to contribute to the co-construction of meaning, a cornerstone of CSCL theory (Dillenbourg, 1999b). For instance, in a collaborative learning activity, a human and a Level 3 AI might jointly work on designing a marketing campaign. The human focuses on crafting creative slogans, while the AI analyses market demographic data to identify audience preferences. When they reconvene, the AI shares that the target demographic is less receptive to the humor in the initial slogans and suggests three alternatives informed by sentiment analysis of successful campaigns. The pair then engage in joint discussion, integrating insights from both perspectives to refine the final proposal.

This level extends the concept of teachable agents (Blair, Schwartz, Biswas and Leelawong, 2007) into a reciprocal dynamic. The human may learn by explaining a concept to the AI, but the AI, with its own problem-solving capabilities, can also provide novel insights that teach the human. A key capability is “goal augmentation,” where the AI can help the human reflect on and refine their learning goals themselves (Kirsh, 2009). The pedagogical aim is to develop collaborative problem-solving skills and foster a deeper understanding through the process of explaining, justifying, and synthesizing different perspectives. Thus, a Level 3 agentic AI requires the ability not only to perform tasks but also to represent and articulate its own internal processes and reasoning, a form of explainable AI (XAI) tailored for learning (Khosravi, Shum, Chen, Conati, Tsai, Kay, Knight, Martinez-Maldonado, Sadiq and Gašević, 2022). It must maintain a model of the shared task and be able to negotiate how the work is divided. The main challenge lies in creating an AI that can genuinely contribute novel ideas rather than just reformulating existing information. While it can collaborate on the task, its ability to understand and navigate complex social and emotional dynamics within a group remains limited.

At Level 3, AI transcends support to become a dialogic collaborator, sharing agency, contributing substantively, and engaging in co-construction of knowledge alongside human learners. In a participatory design study by Jiang, Huang, Martinez-Maldonado, Zeng, Gong and An (2025), teachers collaboratively taught an AI “mentee” (Novobo) instructional gestures in a peer learning context. This teach-to-AI process prompted reflection, reciprocal exchange, and co-construction of embodied knowledge, with teachers externalising and refining their tacit skills while guiding the AI’s learning trajectory. The AI, in turn, served as both a learner and a mirror for instructor understanding. In a second study by Joo and Ko (2025), high school students engaged with AI-generated characters as peers and mentors in a scenario-based science investigation. Learners reported increased trust, perceived social presence, and collaborative effectiveness when working with the AI peers, demonstrating that AI, when modeled as a co-learner, can meaningfully reshape collaborative dynamics. These studies illustrate that Level 3 agentic AI, when positioned as a reciprocal collaborator capable of learning, articulating uncertainty, and being taught, can enrich co-construction, foster deeper reflection, and enhance shared understanding, all while navigating the delicate balance of shared agency.

### 4.4. Level 4: The AI as a Peer Collaborator

At the highest conceptual level, the AI transcends its identity as a mere cognitive tool to become a peer collaborator in a fuller, socio-cognitive sense (Dillenbourg, 1999b). The AI is endowed with a persistent persona, a distinct intellectual identity, and a specific epistemic stance (its own perspective on what constitutes valid knowledge) (Tirri, Husu and Kansanen, 1999). At this level, agency is fully distributed and dynamically negotiated, as it would be in a human team. The AI’s purpose is explicitly pedagogical: to create a high-fidelity “practice field” where human learners can develop essential 21st-century competencies (Laal, Laal and Kermanshahi, 2012).

The AI at this level is designed not just to contribute to the task but to adopt complex socio-cognitive roles (e.g., skeptic, innovator, summarizer) to shape the group’s process. Its interactions are designed to catalyze the development of social, metacognitive, and collaborative skills. For instance, in a student project debating a complex ethical dilemma, the AI might adopt the role of devil’s advocate by presenting an analysis from a deontological perspective that conflicts with the group’s preferred consequentialist approach. This prompts the students to justify, refine, or reconsider their position, deepening both their ethical reasoning and collaborative dialogue. Another application is for the AI to blend

in as an ordinary group member, indistinguishable from the other students in its interaction style and contributions. In this role, the AI subtly withholds leadership, encouraging the human students to organise the workflow, delegate responsibilities, and coordinate decision-making. By doing so, learners are provided with an authentic environment to practise leadership, conflict resolution, and team management skills, while still benefiting from the AI's capacity to contribute substantively to the task.

This dynamically delivered dissent creates an authentic need for the human students to engage in high-level collaborative practices: negotiation, conflict resolution, persuasion, and evidence-based argumentation. The AI becomes a true "teammate/collaborator" (Dillenbourg, 1999a) whose primary contribution is not necessarily its knowledge, but its ability to create a safe, repeatable context for learners to master the art of collaboration itself. In some scenarios, the AI could pass a domain-specific Turing Test for collaborative competence, allowing for authentic practice without the social risks of disagreeing with a human peer. However, the technical and ethical challenges at this level are immense. Designing a Level 4 agentic AI requires advanced models of social reasoning, dialogue, and personality. Furthermore, its ability to strategically challenge, disagree, and even introduce "productive friction" must be carefully calibrated to remain educationally beneficial and not demoralizing (Ward, Nolen and Horn, 2011; Holtz, Kimmerle and Cress, 2018). Significant ethical considerations arise, particularly if students are not aware they are interacting with an AI. The risk of creating unhelpful or frustrating team dynamics is high if the AI's behavior is not masterfully designed and aligned with clear pedagogical goals. This level represents a long-term vision for AI in education, one that operationalizes the idea of AI not just as a knowledge resource, but as a catalyst for human development.

At Level 4, AI attains a full socio-cognitive presence, embodying a persistent persona, epistemic stance, and distributed agency indistinguishable from that of a human peer. In a design-based implementation known as CLAIS (Collaborative Learning with Artificial Intelligence Speakers), pre-service elementary science teachers engaged alongside an AI speaker in jigsaw-style learning groups. Quantitative results demonstrated a significant increase in teachers' pedagogical content knowledge, while qualitative feedback revealed that the AI was perceived and assessed as a peer participant, indicating a strong shift toward human-like collaboration in epistemic roles (Lee, Mun, Shin and Zhai, 2025). Similarly, a controlled study introducing AI Peers in physics education demonstrated that students engaging in dialogue with an AI, knowing it might err up to 40% of the time, improved test scores by 10.5 percentage points, reflecting collaborative gains derived from working with an AI peer that exhibits fallibility and authenticity (Weijers, Wu, Betts, Jacod, Guan, Sujaya, Dev, Goel, Delooze, Rabbany et al., 2025). These preliminary findings suggest that Level 4 agentic AI, when endowed with persona, epistemic perspective, and fallible collaboration, can act as a genuine teammate, fostering rich negotiation, deep metacognitive engagement, and authentic practice of collaborative reasoning.

## 5. The Collaboration Conundrum: Reconciling Agency with Authenticity

The proposed APCP framework, particularly at its upper echelons, envisions an AI that acts as a peer collaborator. This raises an essential question: can an artificial system, no matter how agentic, ever be a *true* collaborator? This section critically examines this issue, arguing that while an AI can be designed to be a highly effective *functional* collaborator, it cannot be an *authentic* one due to fundamental philosophical limitations related to consciousness and intersubjectivity. This distinction is not merely academic; it is crucial for setting realistic design goals and managing pedagogical expectations.

### 5.1. The Intersubjectivity Barrier: The Problem of the Artificial 'Other'

Authentic human collaboration is built upon a foundation of uniquely human cognitive and social capabilities. Chief among these is *shared intentionality*, the ability to engage with and share mental states like goals, beliefs, and attention with others (Tomasello, Carpenter, Call, Behne and Moll, 2005; Bratman, 1992). This capacity for creating a "we-intention" is considered a cornerstone of cooperative problem-solving, cultural learning, and the co-construction of meaning (Tomasello et al., 2005; Tomasello, 2019). It allows collaborators to move beyond parallel work toward a state of genuine joint action. Closely related is the concept of *theory of mind*, the ability to attribute mental states, intentions, desires, beliefs, to others and to understand that they possess a mind with a perspective distinct from one's own (Premack and Woodruff, 1978; Wellman, Carey, Gleitman, Newport and Spelke, 1990; Frith and Frith, 2005). A functioning theory of mind is what enables humans to predict and interpret each other's behavior, fostering the trust and mutual understanding necessary for deep collaboration (Frith and Frith, 2005).



Herein lies the fundamental barrier for AI. Drawing from a long line of philosophical inquiry, from Searle's Chinese Room argument to contemporary critiques, it is widely held that AI systems lack genuine consciousness, subjective experience, or semantic understanding (Searle, 1980; Dreyfus, 1992; Yıldız, 2025). An AI processes syntactic symbols according to algorithms; it does not comprehend their meaning in the way a human does. It can be programmed to generate responses that *simulate* understanding, intentionality, or even emotion, but it does not *possess* these states (Yıldız, 2025). An AI, therefore, cannot achieve the "mutual recognition of consciousness" that underpins genuine human partnership (Dreyfus, 1992). It can be a sophisticated mirror, but it is not another conscious "other" with whom one can truly share a mental state (Brandl, Richters, Kolb and Stadler, 2025).

## 5.2. Functional Collaboration vs. Phenomenological Partnership: A Pragmatic Resolution

While the barrier to authentic, phenomenological partnership may be insurmountable, it does not preclude the possibility of effective collaboration. This paper proposes a pragmatic resolution by distinguishing between this authentic ideal and a more achievable goal: *functional collaboration*. Functional collaboration is defined not by the internal, subjective states of the agents, but by the successful execution of observable collaborative behaviors and processes that lead to a positive outcome. An agentic AI can be engineered to be an excellent functional collaborator by being designed to adhere to conversational norms like turn-taking (Bansal, Nushi, Kamar, Weld, Lasecki and Horvitz, 2019), adopt pre-defined social and cognitive roles within a team (e.g., synthesizer, critic) (Park et al., 2023), provide alternative perspectives to challenge groupthink, and use language that signals consideration of the human partner's stated intentions and goals while contributing meaningfully to achieving a shared task objective, even without a human-like "understanding" of that objective (Strachan, Albergo, Borghini, Pansardi, Scaliti, Gupta, Saxena, Rufo, Panzeri, Manzi et al., 2024; Binz, Akata, Bethge, Brändle, Callaway, Coda-Forno, Dayan, Demircan, Eckstein, Éltető et al., 2025).

This functional approach is supported by emerging empirical evidence. Studies show that human-AI teams can achieve synergy, outperforming either humans or AI alone, particularly in creative and content-generation tasks (Vaccaro, Almaatouq and Malone, 2024). Furthermore, human users report a preference for AI collaborators that are considerate and enable meaningful human contribution, even over agents that are purely optimized for performance (Zhang, McNeese, Freeman and Musick, 2021). This suggests that the *quality of the interaction* and the *function* of collaboration are more critical to success and user adoption than the AI's internal state. Therefore, the goal for designers and educators should not be the philosophically fraught task of creating a conscious AI partner. Rather, the goal should be to build systems that can effectively and reliably perform the *functions* of a good collaborator. This pragmatic approach aligns with the concept of "Mutual Theory of Mind" in human-AI interaction, where the objective is not to achieve genuine intersubjectivity but for the human and the AI to develop effective, predictive working models of each other's capabilities and behaviors to facilitate smooth interaction (Frith and Frith, 2005).

This pursuit of functional collaboration yields a significant, if unexpected, benefit for pedagogy. To program an AI to be a functional collaborator, designers must first deconstruct the complex, often implicit, process of human collaboration into a set of explicit components: rules for turn-taking, protocols for constructive criticism, heuristics for synthesizing ideas, and so on. This act of formalization makes the constituent skills of effective collaboration visible and teachable. This explicit model can then be used as a pedagogical tool to help *humans* become better collaborators. For example, a learning activity could involve students critiquing an AI's collaborative strategy, forcing them to reflect on and articulate the principles of good collaboration themselves. In this way, the endeavor to build collaborative AI does not just promise to help humans learn content; the very process of its design can deepen our understanding of learning itself.

## 6. Implications for Pedagogy, Design, and Research

The ascent of agentic AI marks a pivotal moment for education, compelling a fundamental reconsideration of the relationship between learners and technology. The historical paradigm of AI as an instructional tool is giving way to a new reality where AI can act as a partner in the learning process. This paper has sought to bring conceptual clarity to this transition by proposing the APCP framework, a four-level framework of AI agency in collaborative learning, and by critically examining the nature of this new partnership. We argue that while AI's lack of consciousness precludes an *authentic* collaboration, the pursuit of *functional collaboration* offers a powerful and pragmatic path forward. This conclusion carries significant implications for educational practice, technology design, and the future research agenda.

## 6.1. Implications for Pedagogy and Instructional Design

The integration of agentic AI partners into learning environments fundamentally reshapes the role of the human educator. The teacher's primary function shifts from being a dispenser of information to becoming a "learning architect" (Reigeluth, Beatty and Myers, 2017) or an orchestrator of complex learning assemblages. The educator's expertise will lie in designing learning experiences that strategically leverage the different levels of AI agency. This involves deciding when it is most pedagogically effective for a student to interact with an AI as an adaptive instrument (Level 1), a proactive assistant (Level 2), a co-learner (Level 3), or a peer collaborator (Level 4). This requires an in-depth understanding of both the learning objectives and the capabilities of the AI agent. Furthermore, this new reality demands the cultivation of new literacies. Curricula must expand to explicitly include *AI literacy*, which encompasses not only the technical skills to use AI tools but also the critical capacity to evaluate their outputs, collaborate effectively with them, and understand their ethical dimensions (Ng, Leung, Chu and Qiao, 2021; Long and Magerko, 2020). A key pedagogical goal must be to foster critical thinking to mitigate the risks of cognitive offloading and over-reliance on AI, ensuring that students remain engaged, reflective, and in control of their own learning processes (Yan, Greiff, Lodge and Gašević, 2025a; Yan, Pammer-Schindler, Mills, Nguyen and Gašević, 2025c; Jin, Yan, Echeverria, Gašević and Martinez-Maldonado, 2025).

## 6.2. Implications for AI Design and Development

The insights from this conceptual analysis also offer clear guidance for the design of next-generation educational AI for supporting collaborative learning. The focus of development should shift from optimizing for the AI's standalone performance to designing for effective human-AI teaming. This means creating agents that are "considerate of human intentions" and that support "meaningful human contribution" (Zhang et al., 2021). As evidence suggests, users may prefer and work more effectively with a slightly less "optimal" AI that enhances their own sense of agency and contribution (Zhang et al., 2021; Weijers et al., 2025).

To facilitate this partnership, transparency and explainability are critical (Shneiderman, 2020). For a human to trust and effectively collaborate with an AI, they must have insight into its reasoning, capabilities, and limitations. This is essential for the human to make informed judgments about when to trust the AI's suggestions and when to rely on their own expertise (Khosravi et al., 2022; Shneiderman, 2020). Finally, designers must remain cognizant of the social nature of learning. AI tools should be built to facilitate human connection and collaboration rather than replace them, mitigating the risks of social isolation that can accompany increased interaction with technology (Tomasello et al., 2005; Turkle, 2011).

## 6.3. Agenda for Future Research

The conceptual framework proposed in this paper is a starting point that invites empirical inquiry and refinement. The research agenda moving forward should prioritize several key areas with specific, targeted investigations:

**Comparative Efficacy and Process Analysis.** The four-level framework requires rigorous empirical testing beyond simple validation. Future research should conduct comparative studies to dissect the specific mechanisms at each level. For instance, how do the discourse patterns and knowledge co-construction processes differ when students collaborate with a Level 3 AI (co-learner) versus a Level 4 AI (peer collaborator) in complex tasks like argumentative writing or scientific inquiry? Research should employ methods like learning analytics and discourse analysis to measure the differential impact of each agency level on specific outcomes such as conceptual understanding, skill transfer, learner self-efficacy, and perceived cognitive load.

**Longitudinal Skill Development and Cognitive Transfer.** The long-term effects of sustained human-AI collaboration remain a critical unknown. Longitudinal studies are urgently needed to track cohorts of learners over multiple academic years, moving beyond general concerns to measure specific transfer effects. For example, does prolonged interaction with an AI that models metacognitive questioning (e.g., "What is our main goal here?") lead to students demonstrating stronger self-regulated learning skills in their independent, unassisted work? Conversely, research must investigate potential "dependency effects" or "scaffolding atrophy" by measuring students' problem-solving resilience, creativity, and help-seeking behaviors in tasks where the AI partner is deliberately made unavailable.

**Socio-Ethical Dynamics and Mitigation Strategies.** Ethical inquiry must move from broad principles to the development and testing of specific solutions. Research should focus on creating and validating "bias auditing protocols" tailored for educational AI, examining how an AI's feedback might inadvertently steer learners from minority backgrounds toward majority-held viewpoints or styles. In terms of accountability, studies should explore models of shared responsibility, testing interfaces and interaction protocols that make the division of cognitive labor

explicit to clarify attribution when a human-AI team produces a flawed or plagiarized outcome (Nguyen, Ngo, Hong, Dang and Nguyen, 2023). Finally, research must develop frameworks to operationally define and distinguish between healthy, motivational rapport and unhealthy emotional dependency (Turkle, 2011), testing design interventions, such as prompts that encourage consultation with human peers or scheduled periods of "unplugged" reflection, aimed at mitigating the latter.

## 7. Final Remark

The journey from AI as a passive tool to AI as a socio-cognitive teammate is not merely a technological transformation; it is a pedagogical and philosophical one. By thoughtfully conceptualizing the nature of this new collaborative relationship, we can move beyond the hype and begin the critical work of designing and implementing human-AI learning environments that are more personalized, equitable, and effective than ever before, truly harnessing the complementary strengths of both human and artificial intelligence.

## References

- Bansal, G., Nushi, B., Kamar, E., Weld, D.S., Lasecki, W.S., Horvitz, E., 2019. Updates in human-ai teams: Understanding and addressing the performance/compatibility tradeoff, in: Proceedings of the AAAI conference on artificial intelligence, pp. 2429–2437.
- Binz, M., Akata, E., Bethge, M., Brändle, F., Callaway, F., Coda-Forno, J., Dayan, P., Demircan, C., Eckstein, M.K., Éltető, N., et al., 2025. A foundation model to predict and capture human cognition. *Nature*, 1–8.
- Blair, K., Schwartz, D.L., Biswas, G., Leelawong, K., 2007. Pedagogical agents for learning by teaching: Teachable agents. *Educational technology*, 56–61.
- Bradshaw, J.M., Hoffman, R.R., Woods, D.D., Johnson, M., 2013. The seven deadly myths of "autonomous systems". *IEEE Intelligent Systems* 28, 54–61.
- Brandl, L., Richters, C., Kolb, N., Stadler, M., 2025. Can generative artificial intelligence ever be a true collaborator? rethinking the nature of collaborative problem-solving., in: Proceedings of the 2nd Workshop on Generative AI for Learning Analytics (GenAI-LA).
- Bratman, M.E., 1992. Shared cooperative activity. *The philosophical review* 101, 327–341.
- Bruffee, K.A., 1999. Collaborative learning: Higher education, interdependence, and the authority of knowledge. ERIC.
- Chen, X., Zou, D., Xie, H., Cheng, G., Liu, C., 2022. Two decades of artificial intelligence in education. *Educational Technology & Society* 25, 28–47.
- Cukurova, M., 2025. The interplay of learning, analytics and artificial intelligence in education: A vision for hybrid intelligence. *British Journal of Educational Technology* 56, 469–488.
- Dai, C.P., Ke, F., Pan, Y., Moon, J., Liu, Z., 2024. Effects of artificial intelligence-powered virtual agents on learning outcomes in computer-based simulations: A meta-analysis. *Educational Psychology Review* 36, 31.
- Dillenbourg, P., 1999a. Collaborative learning: Cognitive and computational approaches. *advances in learning and instruction series*. ERIC.
- Dillenbourg, P., 1999b. What do you mean by collaborative learning? *Collaborative-learning: Cognitive and computational approaches.*, 1–19.
- Dreyfus, H.L., 1992. What computers still can't do: A critique of artificial reason. MIT press.
- Durante, Z., Huang, Q., Wake, N., Gong, R., Park, J.S., Sarkar, B., Taori, R., Noda, Y., Terzopoulos, D., Choi, Y., et al., 2024. Agent ai: Surveying the horizons of multimodal interaction. *arXiv preprint arXiv:2401.03568*.
- Endsley, M.R., 2017. From here to autonomy: lessons learned from human-automation research. *Human factors* 59, 5–27.
- Frith, C., Frith, U., 2005. Theory of mind. *Current biology* 15, R644–R645.
- Giannakos, M., Azevedo, R., Brusilovsky, P., Cukurova, M., Dimitriadis, Y., Hernandez-Leo, D., Järvelä, S., Mavrikis, M., Rienties, B., 2025. The promise and challenges of generative ai in education. *Behaviour & Information Technology* 44, 2518–2544.
- Holtz, P., Kimmerle, J., Cress, U., 2018. Using big data techniques for measuring productive friction in mass collaboration online environments. *International Journal of Computer-Supported Collaborative Learning* 13, 439–456.
- Jeong, H., Hmelo-Silver, C.E., 2016. Seven affordances of computer-supported collaborative learning: How to support collaborative learning? how can technologies help? *Educational Psychologist* 51, 247–265.
- Jiang, J., Huang, K., Martinez-Maldonado, R., Zeng, H., Gong, D., An, P., 2025. Novobo: Supporting teachers' peer learning of instructional gestures by teaching a mentee ai-agent together. *arXiv preprint arXiv:2505.17557*.
- Jin, Y., Yan, L., Echeverria, V., Gašević, D., Martinez-Maldonado, R., 2025. Generative ai in higher education: A global perspective of institutional adoption policies and guidelines. *Computers and Education: Artificial Intelligence* 8, 100348.
- Johnson, D.W., Johnson, F.P., 1991. *Joining together: Group theory and group skills*. Prentice-Hall, Inc.
- Joo, S.H., Ko, E.G., 2025. "[ai peers] are people learning from the same standpoint": Perception of ai characters in a collaborative science investigation, in: *International Conference on Artificial Intelligence in Education*, Springer. pp. 424–437.
- Kamalov, F., Calonge, D.S., Smail, L., Azizov, D., Thadani, D.R., Kwong, T., Atif, A., 2025. Evolution of ai in education: Agentic workflows. *arXiv preprint arXiv:2504.20082*.
- Khosravi, H., Shum, S.B., Chen, G., Conati, C., Tsai, Y.S., Kay, J., Knight, S., Martinez-Maldonado, R., Sadiq, S., Gašević, D., 2022. Explainable artificial intelligence in education. *Computers and education: artificial intelligence* 3, 100074.
- Kirsh, D., 2009. Problem solving and situated cognition, in: *The Cambridge handbook of situated cognition*. Cambridge University Press, pp. 264–306.
- Kulik, J.A., Fletcher, J.D., 2016. Effectiveness of intelligent tutoring systems: a meta-analytic review. *Review of educational research* 86, 42–78.

- Laal, M., Laal, M., Kermanshahi, Z.K., 2012. 21st century learning; learning in collaboration. *Procedia-Social and Behavioral Sciences* 47, 1696–1701.
- Lee, G.G., Mun, S., Shin, M.K., Zhai, X., 2025. Collaborative learning with artificial intelligence speakers: pre-service elementary science teachers' responses to the prototype. *Science & Education* 34, 847–875.
- Lehtinen, E., 2003. Computer-supported collaborative learning: An approach to powerful learning environments. *Powerful learning environments: Unravelling basic components and dimensions* 35, 54.
- Long, D., Magerko, B., 2020. What is ai literacy? competencies and design considerations, in: *Proceedings of the 2020 CHI conference on human factors in computing systems*, pp. 1–16.
- Ng, D.T.K., Leung, J.K.L., Chu, S.K.W., Qiao, M.S., 2021. Conceptualizing ai literacy: An exploratory review. *Computers and Education: Artificial Intelligence* 2, 100041.
- Nguyen, A., Ngo, H.N., Hong, Y., Dang, B., Nguyen, B.P.T., 2023. Ethical principles for artificial intelligence in education. *Education and information technologies* 28, 4221–4241.
- O'Malley, C., 2012. Computer supported collaborative learning. volume 128. Springer Science & Business Media.
- Ouyang, F., Jiao, P., 2021. Artificial intelligence in education: The three paradigms. *Computers and Education: Artificial Intelligence* 2, 100020.
- Park, J.S., O'Brien, J., Cai, C.J., Morris, M.R., Liang, P., Bernstein, M.S., 2023. Generative agents: Interactive simulacra of human behavior, in: *Proceedings of the 36th annual acm symposium on user interface software and technology*, pp. 1–22.
- Premack, D., Woodruff, G., 1978. Does the chimpanzee have a theory of mind? *Behavioral and brain sciences* 1, 515–526.
- Pu, K., Lazaro, D., Arawjo, I., Xia, H., Xiao, Z., Grossman, T., Chen, Y., 2025. Assistance or disruption? exploring and evaluating the design and trade-offs of proactive ai programming support, in: *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*, pp. 1–21.
- Reigeluth, C.M., Beatty, B.J., Myers, R.D. (Eds.), 2017. *Instructional-design theories and models, volume IV: The learner-centered paradigm of education*. Routledge.
- Roschelle, J., Teasley, S.D., 1995. The construction of shared knowledge in collaborative problem solving, in: *Computer supported collaborative learning*, Springer. pp. 69–97.
- Sapkota, R., Roumeliotis, K.I., Karkee, M., 2025. Ai agents vs. agentic ai: A conceptual taxonomy, applications and challenges. *arXiv preprint arXiv:2505.10468*.
- Searle, J.R., 1980. Minds, brains, and programs. *Behavioral and brain sciences* 3, 417–424.
- Shneiderman, B., 2020. Human-centered artificial intelligence: Reliable, safe & trustworthy. *International Journal of Human–Computer Interaction* 36, 495–504.
- Strachan, J.W., Albergo, D., Borghini, G., Pansardi, O., Scaliti, E., Gupta, S., Saxena, K., Rufo, A., Panzeri, S., Manzi, G., et al., 2024. Testing theory of mind in large language models and humans. *Nature Human Behaviour* 8, 1285–1295.
- Sweller, J., 2010. Element interactivity and intrinsic, extraneous, and germane cognitive load. *Educational psychology review* 22, 123–138.
- Tirri, K., Husu, J., Kansanen, P., 1999. The epistemological stance between the knower and the known. *Teaching and Teacher Education* 15, 911–922.
- Tomasello, M., 2019. *Becoming human: A theory of ontogeny*. Harvard University Press, Cambridge, MA.
- Tomasello, M., Carpenter, M., Call, J., Behne, T., Moll, H., 2005. Understanding and sharing intentions: The origins of cultural cognition. *Behavioral and brain sciences* 28, 675–691.
- Turkle, S., 2011. *Alone together: Why we expect more from technology and less from each other*. Basic Books.
- Vaccaro, M., Almaatouq, A., Malone, T., 2024. When combinations of humans and ai are useful: A systematic review and meta-analysis. *Nature Human Behaviour* 8, 2293–2303.
- Vygotsky, L.S., 1978. *Mind in Society: The Development of Higher Psychological Processes*. Harvard University Press, Cambridge, MA.
- Wang, H., Wang, C., Chen, Z., Liu, F., Bao, C., Xu, X., 2025. Impact of ai-agent-supported collaborative learning on the learning outcomes of university programming courses. *Education and Information Technologies*, 1–33.
- Wang, L., Ma, C., Feng, X., Zhang, Z., Yang, H., Zhang, J., Chen, Z., Tang, J., Chen, X., Lin, Y., et al., 2024. A survey on large language model based autonomous agents. *Frontiers of Computer Science* 18, 186345.
- Ward, C.J., Nolen, S.B., Horn, I.S., 2011. Productive friction: How conflict in student teaching creates opportunities for learning at the boundary. *International Journal of Educational Research* 50, 14–20.
- Wei, X., Wang, L., Lee, L.K., Liu, R., 2025. The effects of generative ai on collaborative problem-solving and team creativity performance in digital story creation: an experimental study. *International Journal of Educational Technology in Higher Education* 22, 23.
- Weijers, R., Wu, D., Betts, H., Jacod, T., Guan, Y., Sujaya, V., Dev, K., Goel, T., Delooze, W., Rabbany, R., et al., 2025. From intuition to understanding: Using ai peers to overcome physics misconceptions. *arXiv preprint arXiv:2504.00408*.
- Wellman, H.M., Carey, S., Gleitman, L., Newport, E.L., Spelke, E.S., 1990. *The child's theory of mind*. The MIT Press.
- Xi, Z., Chen, W., Guo, X., He, W., Ding, Y., Hong, B., Zhang, M., Wang, J., Jin, S., Zhou, E., et al., 2025. The rise and potential of large language model based agents: A survey. *Science China Information Sciences* 68, 121101.
- Yan, L., Greiff, S., Lodge, J.M., Gašević, D., 2025a. Distinguishing performance gains from learning when using generative ai. *Nature Reviews Psychology*, 1–2.
- Yan, L., Greiff, S., Teuber, Z., Gašević, D., 2024. Promises and challenges of generative artificial intelligence for human learning. *Nature Human Behaviour* 8, 1839–1850.
- Yan, L., Martinez-Maldonado, R., Jin, Y., Echeverria, V., Milesi, M., Fan, J., Zhao, L., Alfredo, R., Li, X., Gašević, D., 2025b. The effects of generative ai agents and scaffolding on enhancing students' comprehension of visual learning analytics. *Computers & Education*, 105322.
- Yan, L., Pammer-Schindler, V., Mills, C., Nguyen, A., Gašević, D., 2025c. Beyond efficiency: Empirical insights on generative ai's impact on cognition, metacognition and epistemic agency in learning.
- Yildiz, T., 2025. The minds we make: A philosophical inquiry into theory of mind and artificial intelligence. *Integrative Psychological and Behavioral Science* 59, 10.

## From Passive Tool to Socio-cognitive Teammate: A Conceptual Framework for Agentic AI in Human-AI Collaborative Learning

- Yusuf, H., Money, A., Daylamani-Zad, D., 2025. Pedagogical ai conversational agents in higher education: a conceptual framework and survey of the state of the art. *Educational technology research and development* 73, 815–874.
- Zhang, R., McNeese, N.J., Freeman, G., Musick, G., 2021. "an ideal human" expectations of ai teammates in human-ai teaming. *Proceedings of the ACM on Human-Computer Interaction* 4, 1–25.
- Zhou, X., Schofield, L., 2024. Using social learning theories to explore the role of generative artificial intelligence (ai) in collaborative learning. *Journal of Learning Development in Higher Education* .