

Review

Beyond Automation: Pedagogical Strategies for Meaningful Human-AI Collaboration in the Classroom

Dongxing Yu 

Institute for Sustainable Development in Education, School of Education, Shanghai Sanda University, Shanghai 201209, China; yudongxing@sandau.edu.cn

Received: 30 August 2025; **Revised:** 16 November 2025; **Accepted:** 2 December 2025; **Published:** 4 January 2026

Abstract: This article examines how generative artificial intelligence can be integrated into teaching without reducing learning to automated answer production. Rather than treating AI adoption as a purely technical question, the paper argues that the central pedagogical challenge is how to preserve human judgment, metacognitive effort, and disciplinary understanding when students can now generate plausible outputs in seconds. Building on the original manuscript's distinction between automation and augmentation, the revised version strengthens the argument by anchoring the problem in concrete classroom pain points, especially writing-intensive and feedback-heavy courses where students' use of AI often outpaces institutional policy. It synthesizes recent scholarship on hybrid intelligence, self-regulated learning, teacher-AI collaboration, assessment redesign, and academic integrity, while also identifying limitations in current frameworks, including weak subject-specific guidance, limited long-term evidence, and insufficient attention to equity and cultural context. In response, the paper clarifies the rationale for the COLLABORATE framework and explains how its ten principles work together to make AI use visible, bounded, and pedagogically productive. The framework is presented as a conceptual, evidence-informed design model rather than as an empirically validated intervention. The paper concludes by outlining practical implementation pathways, ethical safeguards, and a concrete research agenda for testing which forms of human-AI collaboration best support student learning, process transparency, and foundational skill retention.

Keywords: Human-AI Collaboration; Generative AI; Pedagogical Strategies; AI Literacy; Metacognition; Teacher-AI Partnership; Educational Technology

1. Introduction

Students can now produce credible academic prose, code, and explanations in seconds, yet many classrooms still assess learning as if visible final products reliably reflect students' own thinking. This contradiction is not marginal; it sits at the center of the current academic integrity crisis. Since the public release of ChatGPT in late 2022, generative AI has moved from peripheral educational technology to a direct participant in cognitive labor, forcing educators to reconsider authorship, assessment, and skill formation [1,2].

The pain points become especially tangible in writing-intensive instruction. In a first-year academic writing course, for example, students may use AI to brainstorm, outline, paraphrase, edit, or fully draft assignments long before institutions have clarified what constitutes acceptable support. UNESCO [3] reported that fewer than 10% of schools and universities worldwide had formal guidance on AI, while more recent surveys suggest that student adoption has outpaced policy responses and teacher preparation [4,5]. The result is a shadow curriculum in which students experiment with powerful tools under conditions of low transparency and uneven guidance.

Existing responses have not solved this problem. Prohibition-based approaches are difficult to enforce, AI-detection tools remain unreliable, and generic policy statements rarely explain how teachers should redesign tasks so that AI supports rather than replaces disciplinary thinking. At the same time, scholarship has raised legitimate concerns about cognitive offloading, superficial revision, bias reproduction, and widening inequities in access and outcomes [6,7]. What is still missing is not another tool guide, but a pedagogically grounded framework that shows how human-AI collaboration can be structured differently across authentic classroom contexts.

This paper addresses that gap by arguing that the key question is not whether AI belongs in classrooms, but under what conditions its use strengthens rather than weakens human learning. Accordingly, the article moves from the problem of automation to the possibility of augmentation, then derives the COLLABORATE framework as a design response to four linked problems: invisible process, weak metacognitive engagement, assessment misalignment, and insufficient teacher support. The paper’s concrete objectives are to clarify the framework’s theoretical foundations, explain why existing approaches remain insufficient, specify how the ten principles relate to one another, and outline verifiable research questions for future empirical testing.

2. Theoretical Foundations

2.1. From Automation to Augmentation

The distinction between automation and augmentation serves as the primary epistemological pivot point for understanding the role of artificial intelligence in educational contexts. Automation, in its strictest sense, involves the deployment of technology to execute tasks previously performed by humans, with the explicit aim of increasing utilitarian efficiency or reducing labor costs. In an educational setting, this manifests as systems that grade essays without human oversight, generate lesson plans without pedagogical customization, or produce student outputs that bypass the cognitive struggle of creation. While automation offers administrative convenience, it carries the inherent risk of “de-skilling,” wherein the learner—or indeed the educator—cedes critical competencies to the machine, resulting in a degradation of human agency and expertise.

In diametric contrast, augmentation—often referred to in literature as Intelligence Amplification (IA)—positions the artificial intelligence not as a substitute for human cognition, but as a scaffolding mechanism that extends human capabilities. This paradigm posits that the symbiotic pairing of human judgment and algorithmic processing can achieve outcomes superior to either entity operating in isolation [8]. Augmentation does not seek to remove the human from the loop; rather, it seeks to elevate the human’s entry point into the problem. By offloading routine information retrieval or pattern recognition to the AI, the learner is liberated to focus on higher-order cognitive functions such as synthesis, ethical evaluation, and contextual application. **Table 1** summarizes the contrast between automation and augmentation in AI-enhanced education.

Table 1. Comparison of Automation and Augmentation Approaches in AI-Enhanced Education.

Dimension	Automation Approach	Augmentation Approach
Primary Goal	Replace human effort (Efficiency focus)	Enhance human capabilities (Efficacy focus)
Student Role	Passive consumer of AI outputs	Active collaborator and evaluator of AI
Teacher Role	Reduced or replaced by AI systems	Orchestrator of human-AI learning ecosystems
Skill Development	Atrophy through disuse (“De-skilling”)	Enhanced through scaffolded practice
Critical Thinking	Bypassed; cognitive shortcuts encouraged	Developed through critique of AI outputs
Learning Outcome	Task completion and speed	Deep understanding and intellectual growth

Note: Synthesis adapted from Wilson and Daugherty [8] and Mollick and Mollick [9].

Effective augmentation relies on the concept of “Hybrid Intelligence,” defined as the ability of a socio-technical system to solve complex problems by leveraging the complementary strengths of human and artificial intelligence [10]. AI excels at processing vast datasets, ensuring syntactical correctness, and generating divergent possibilities at scale. Humans, conversely, excel at semantic understanding, normative judgment, empathy, and the application of knowledge to novel, ambiguous contexts. In this view, the goal of education is not to compete with AI in domains of calculation and retrieval, but to build the metacognitive capacity to manage, direct, and audit these powerful systems.

2.2. Metacognition and Self-Regulated Learning

Metacognition—the awareness and regulation of one’s own cognitive processes—emerges as the sine qua non of effective human-AI collaboration. In traditional learning environments, metacognition involves monitoring one’s understanding of the material. In an AI-integrated environment, this expands to include “Algorithm-Aware Metacognition”: the ability to monitor the quality of the tool’s output, calibrate one’s trust in the system, and discern when to accept, reject, or modify algorithmic suggestions. Research demonstrates that students who possess high levels of metacognitive awareness are more likely to use AI tools productively, engaging in “epistemic checking” rather than passive acceptance [6].

Self-regulated learning (SRL) provides the operational framework for this engagement. SRL involves three cyclical phases: forethought (planning), performance (execution), and self-reflection. Without rigorous self-regulation, students are susceptible to “detrimental cognitive offloading,” a phenomenon where the learner outsources the thinking process entirely to the AI, resulting in no “cognitive residual”—that is, no lasting change in the learner’s long-term memory or skill set.

However, a nuanced understanding distinguishes this from “beneficial cognitive offloading.” Just as a mathematician uses a calculator to offload arithmetic computation to focus on complex proofs, a student might use AI to offload the structural organization of notes to focus on conceptual synthesis. The critical variable is intentionality. When properly scaffolded, AI tools can actually enhance SRL by acting as an external monitor, providing metacognitive prompts (e.g., “Have you considered an alternative perspective?”) that encourage students to reflect on their learning processes rather than simply rushing to a conclusion [9]. Thus, the pedagogical imperative is to design interactions where AI serves as a mirror for the student’s thinking, rather than a mask for their lack of it.

Taken together, these perspectives imply that effective classroom AI use cannot be evaluated only by speed, convenience, or surface-quality output. It must be judged by whether students still engage in planning, monitoring, evaluating, and transferring across tasks. This is why the move from automation to augmentation is not merely semantic: it identifies whether AI is substituting for core learning processes or supporting them in bounded, inspectable ways.

3. Literature Review: Human-AI Collaboration in Education

3.1. Teacher-AI Collaboration Models

Recent scholarship has identified multiple models for teacher-AI collaboration, decisively shifting the narrative from one of inevitable displacement to one of strategic partnership. While early discourse was dominated by fears of AI rendering educators obsolete, contemporary frameworks emphasize the “Teacher-in-the-Loop” approach [11]. In this paradigm, the educator acts not merely as a user of technology, but as the “conductor” of a complex cognitive system, orchestrating the interplay between algorithmic efficiency and human judgment.

This “Teacher-in-the-Loop” model fundamentally redefines the division of labor in the classroom. The AI functions as a force multiplier, adept at handling high-volume, low-complexity tasks that typically drain teacher bandwidth. For example, generative AI can instantaneously produce thirty distinct variations of a practice problem, conduct initial syntax checks on hundreds of student essays, or retrieve specific data points from vast corpora. By automating these routine administrative and logistical burdens, the system frees the teacher to focus on high-complexity, high-impact interactions that machines cannot replicate. These include emotional mentorship, the detection of subtle behavioral nuances indicating struggle, and the provision of ethical guidance during complex discussions.

Furthermore, this partnership allows for a reduction in the “cognitive load” associated with classroom management, enabling teachers to reinvest their energy into “pedagogical reasoning”—the moment-to-moment decisions about how to explain a concept or support a struggling student. The teacher remains the central decision-maker, interpreting AI-generated insights (such as a dashboard flagging a drop in engagement) and deciding whether to intervene with a new activity, a one-on-one conversation, or a class-wide pause. Thus, the goal is not to remove the teacher, but to “up-level” their role from content delivery to learning architecture.

We synthesize emerging practices into six distinct models of collaboration (**Table 2**). These models move along a continuum of AI autonomy, from passive observation to active co-teaching.

Table 2. Synthesized Models of Teacher-AI Collaboration.

Collaboration Model	Description
“One Teach, One Observe”	AI monitors student engagement/data while the teacher delivers instruction.
“One Teach, One Assist”	AI provides real-time tutoring support to individuals while the teacher focuses on the whole class.
Co-teaching in Stations	Students rotate between teacher-led discussion and AI-led practice/drills.
Parallel Teaching	Coordinating AI-facilitated online modules with teacher-led offline reflection.
Differentiated Teaching	AI adapts content difficulty for individual learners; teacher manages the pedagogical strategy.
Team Teaching	Teacher and AI jointly deliver instruction; the teacher prompts the AI live to demonstrate concepts.

Note: Synthesized from emerging frameworks [5,9].

To fully leverage these models, educators must understand the distinct operational dynamics of each. For instance, in the “Team Teaching” model, the AI is often projected on a screen “live,” serving as a third party in the classroom dialogue. A teacher might prompt the AI to “argue against my point” or “summarize this confusion,” effectively modeling intellectual humility and critical interrogation for students. This turns the AI into a pedagogical prop that makes thinking visible. Similarly, the “One Teach, One Assist” model fundamentally alters classroom management; rather than the teacher circulating to answer repetitive clarification questions, the AI handles these “Tier 1” queries, allowing the teacher to target “Tier 2” and “Tier 3” interventions for students with deeper conceptual struggles.

However, these models also introduce new pedagogical demands. The “Differentiated Teaching” approach, while powerful for personalization, carries the risk of algorithmic bias or “tracking,” where the AI might consistently serve lower-level content to struggling students without the teacher’s awareness. Therefore, the “One Teach, One Observe” dynamic becomes crucial: the teacher must constantly audit the AI’s “behavior” just as they would observe a student teacher, ensuring that efficiency does not come at the cost of equity or depth. This requires a shift in teacher identity from “content deliverer” to “learning architect,” responsible for designing the ecosystem in which human and machine intelligence interact.

3.2. Student-AI Interaction Patterns

Research on student-AI interaction has revealed patterns far more complex than simple query-response dynamics. Rather than a linear path from question to answer, students often engage in recursive, non-linear learning cycles when effectively working with AI tools. Ideally, this interaction takes the form of “co-editing” or “co-creation,” a symbiotic process where the student functions as the senior editor and the AI as the junior drafter. In this model, the student provides the strategic intent, structural outline, and tonal constraints, while the AI generates the raw textual material. The student then engages in rigorous critique, factual verification, and stylistic refinement [12]. This iterative loop—prompt, review, refine, re-prompt—mirrors the cognitive processes of expert revision, potentially accelerating the development of higher-order writing skills by allowing students to prototype ideas rapidly.

However, without specific pedagogical guidance, interaction patterns can easily devolve into a passive “paste-and-submit” workflow. In these instances, students treat the AI as an oracle rather than a tool, accepting outputs uncritically and bypassing the cognitive struggle necessary for learning. This “epistemic dependence” poses a significant risk to learning autonomy. Prather et al. [13] found that students explicitly need instruction on how to interact with code-generating models to avoid superficial learning. Their findings suggest that “prompt literacy” is now a prerequisite for academic success, encompassing not just the technical ability to formulate a query, but the metacognitive ability to decompose complex problems into sub-tasks that an AI can effectively handle and the discipline to scrutinize the results.

3.3. AI’s Impact on Critical Thinking

The relationship between AI use and critical thinking remains one of the most contested areas in current educational discourse. Skeptics argue that AI acts as a “cognitive crutch,” leading to the atrophy of fundamental skills like outlining, summarizing, and drafting. If the machine does the synthesizing, the student may never learn to synthesize themselves. Proponents, however, counter that AI acts as a “calculator for writing”—a tool that automates lower-order cognitive tasks (like sentence construction or basic code syntax) to free up mental bandwidth for higher-order thinking (like complex argumentation, logic checking, and creative synthesis).

A study by Gero et al. [14], conducted even before the ubiquitous adoption of ChatGPT, found that AI tools (referred to as “sparks”) could significantly enhance creativity, but only if the system was designed to support human agency rather than overwrite it. The AI served to overcome “blank page syndrome,” offering diverse divergent options that the human user could then converge upon. The danger arises when the AI’s “fluency” is mistaken for “accuracy,” leading students to overestimate the quality of the reasoning.

Consequently, the key moderating factor determining whether AI helps or harms is instructional design. AI tends to undermine critical thinking when the task is merely information retrieval or generic synthesis—tasks the AI performs with deceptive ease. Conversely, AI enhances critical thinking when the task is designed to require evaluation, synthesis of disparate sources, and personal reflection. For example, assignments that ask students to critique AI-generated hallucinations or bias force them to exercise distinctively human judgment, transforming the AI from a source of truth into an object of study.

3.4. Gaps in Existing Frameworks and Rationale for COLLABORATE

Although the literature increasingly favors human-AI partnership over simple replacement, several limitations remain. First, many frameworks operate as tool guides: they explain how to prompt, disclose, or cite AI, but provide limited subject-specific guidance on how to redesign learning sequences. Second, evidence on long-term effectiveness remains thin, particularly with respect to foundational skill retention, revision quality, and transfer. Third, cultural and institutional differences are often underexamined, even though access, trust, language background, and local policy conditions shape how AI is experienced in classrooms. Recent student perception reports and policy analyses likewise show that institutional guidance remains uneven and often reactive [15,16]. Finally, the field still lacks sufficiently explicit design logic showing how assessment, feedback, metacognition, and teacher orchestration must work together rather than as isolated principles [17–20].

The revised manuscript therefore positions COLLABORATE as an evidence-informed conceptual framework rather than as a finished or universally validated model. Its contribution lies in integrating these dimensions into a single pedagogical design logic: AI should be contextualized, bounded, made transparent, and tied to tasks that still require human judgment. In this sense, the framework is intended to connect rather than replace work on self-regulated learning, formative assessment, hybrid intelligence, and AI literacy.

4. The COLLABORATE Framework

Predicated upon a synthesis of extant scholarly research, the COLLABORATE framework is here clarified as a design logic rather than a flat checklist. The ten principles were not selected arbitrarily; they were consolidated from recurring concerns in the literature on hybrid intelligence, self-regulated learning, assessment reform, AI literacy, and teacher orchestration. Together, they address four functional needs: grounding AI in disciplinary context, preserving learner agency, making the process visible, and equipping teachers to supervise and redesign practice. Several principles are interdependent. For example, contextualized learning and balanced integration establish the conditions under which orchestrated interactions and orchestrated feedback can be effective, while reflective practice and assessment redesign make those interactions auditable. **Figure 1** presents the framework and the relationships among its principles.

1. **Contextualized Learning:** It is strictly mandated that the integration of artificial intelligence be firmly and irrevocably situated within authentic, discipline-specific paradigms. This requirement thereby necessitates the eschewing of generic, context-independent applications which frequently result in superficial engagement. Rather, the algorithmic tools must be embedded within the epistemological frameworks unique to each field of study, ensuring that the computational output serves to reinforce, rather than dilute, the specific modes of inquiry and validation characteristic of the discipline.
2. **Orchestrated Interactions:** The active administration of the transitional interfaces between human labor and algorithmic computation is to be executed by the educator. The pedagogue must assume the role of an architect, rigorously defining the temporal sequence of operations—specifically, determining at which junctures artificial intelligence is introduced and at which points it must be withdrawn. This orchestration prevents the technology from dominating the learning trajectory, ensuring instead that it remains a subservient instrument subject to human pedagogical goals.

3. **Learner Agency:** It is requisite that the student retains the position of primary agent within the operational loop, thereby exercising final authority over editorial adjudications. The learner must be positioned not as a passive recipient of generated content, but as the executive decision-maker responsible for the synthesis, verification, and final approval of all outputs. This structure is intended to mitigate the risks associated with “automation bias,” wherein users unthinkingly defer to algorithmic suggestions, thereby eroding independent cognitive function.
4. **Literacy Development:** Instruction regarding artificial intelligence literacy, encompassing both technical and ethical dimensions, must be explicitly delivered rather than presumed. It is insufficient to merely provide access to tools; curricula must systematically dismantle the “black box” nature of these systems. This involves rigorous training in the identification of algorithmic hallucinations, the recognition of embedded biases, and the comprehension of the statistical—rather than sentient—nature of large language models, thereby inoculating students against anthropomorphism.
5. **Assessment Redesign:** Evaluative protocols shall be restructured to accord value to the iterative process of creation, in preference to the static final deliverable [6]. Given the capacity of artificial intelligence to generate high-fidelity final products, the locus of assessment must shift toward the documentation of intellectual evolution. This necessitates the evaluation of draft iterations, prompt histories, and the rationale behind the rejection or acceptance of algorithmic suggestions, effectively rendering the “journey” of learning more significant than the “destination” of the submission.
6. **Balanced Integration:** The maintenance of designated zones wherein artificial intelligence is prohibited is deemed essential for the assurance of foundational skill acquisition. To prevent the atrophy of core cognitive competencies, specific instructional periods must be ring-fenced as “unassisted performance” zones. In these domains, students are compelled to operate without technological scaffolding, ensuring that they possess the requisite mental models and fundamental knowledge necessary to effectively evaluate and guide artificial intelligence in subsequent, augmented tasks.
7. **Orchestrated Feedback:** Utilization of artificial intelligence shall be directed toward the provision of immediate, formative feedback, thus conserving educator resources for high-level conceptual mentorship. While algorithmic systems are capable of providing instantaneous correction on mechanical and syntactical errors at scale, they lack the nuanced understanding required for deep conceptual guidance. Consequently, a division of labor is proposed wherein artificial intelligence manages the lower-order feedback loop, thereby liberating the human educator to address higher-order reasoning, argumentation, and complex structural critique [21].
8. **Reflective Practice:** The documentation of the methods by which artificial intelligence is employed, such as the submission of interaction transcripts, is required of the student. It is not sufficient for the learner to merely utilize the tool; they must also engage in a metacognitive analysis of that utilization. By compelling students to submit logs of their interactions alongside their final work, educators can enforce a transparent audit trail that reveals the extent of algorithmic contribution and the students’ own critical engagement with the machine’s output.
9. **Authentic Tasks:** Assigned tasks must necessitate the application of knowledge to local or personal contexts that are beyond the cognizance of the artificial intelligence. Since large language models are trained on historical, generalized datasets, they inherently lack access to real-time, hyper-local, or phenomenological experiences. Assignments should therefore be designed to exploit this limitation, requiring students to synthesize general theoretical knowledge with specific, embodied realities—such as local community interviews or personal case studies—that the artificial intelligence cannot hallucinate or retrieve.
10. **Teacher Empowerment:** It is incumbent upon the institution to ensure educators are trained to audit and supervise algorithmic recommendations. Teachers must not be reduced to passive implementers of technological policy but must be equipped with the technical competence and administrative authority to challenge, modify, or reject AI-driven interventions. This empowerment ensures that professional pedagogical judgment remains the supreme arbiter in the classroom, preventing the cession of educational sovereignty to automated systems.

Read together, these principles suggest an order of implementation rather than a simple list. Educators typically begin by defining the disciplinary purpose of AI use and the boundaries of acceptable assistance, then design interaction roles, feedback loops, and assessment artifacts that keep student reasoning visible. This dependency

structure also explains why the framework does not add separate principles for every important concern; issues such as equity, cultural responsiveness, and emotional safety should be treated as cross-cutting design criteria applied across contextualization, agency, assessment, and teacher empowerment.

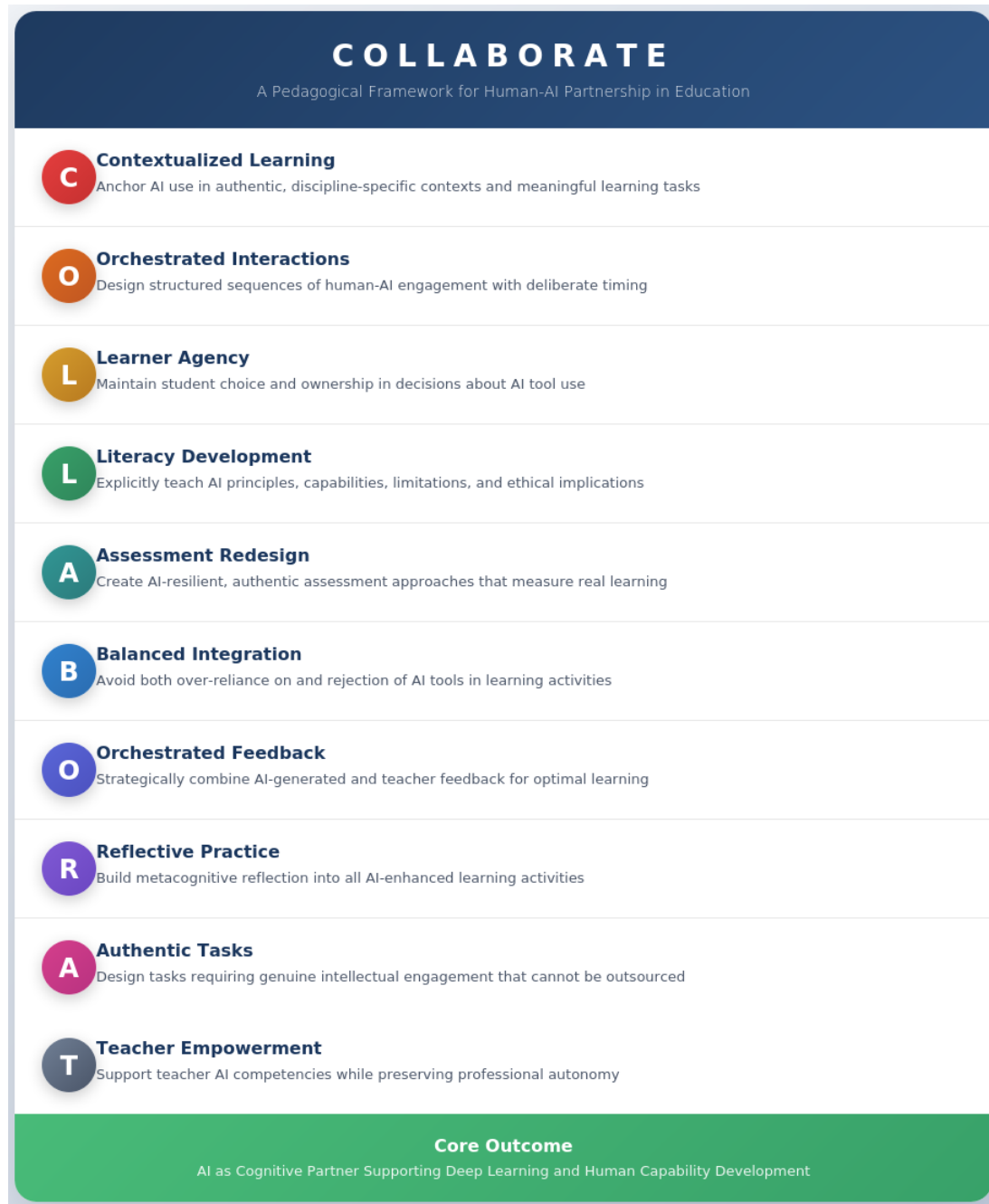


Figure 1. The COLLABORATE Framework for Meaningful Human-AI Collaboration in Education.

Note: The framework emphasizes interconnected principles for meaningful human-AI collaboration.

5. Implementation Strategies

To make the framework more concrete while preserving the original structure of the paper, the revised manuscript emphasizes that implementation must be scenario-specific. A writing course, a physics class, and a history seminar will not use the same collaboration model because the underlying cognitive work differs. The sections below therefore retain the practical orientation of the original manuscript while clarifying that prompt design and scaffolding

are means to preserve disciplinary thinking, not ends in themselves.

5.1. Prompt Engineering as Pedagogy

“Prompt engineering” is frequently mischaracterized in public discourse as a purely technical or utilitarian skill—a “hack” to maximize productivity. However, in a pedagogical context, it functions as a rigorous exercise in linguistic precision, clarity of thought, and structural logic. It represents a new form of literacy that compels the student to externalize their internal cognitive structures. To prompt effectively, a student must first decompose a complex, multifaceted problem into its constituent parts, articulate precise criteria for success, and anticipate potential areas of ambiguity or misunderstanding—cognitive processes that mirror high-level expository writing and computational thinking. Therefore, instruction in prompting is synonymous with instruction in critical thinking. The CLEAR framework [22] provides a robust rubric for this practice, ensuring that students engage in a disciplined, iterative interaction rather than a casual, one-off query. The components of this framework are summarized in Table 3:

Table 3. The CLEAR Framework for Effective Prompt Engineering in Education.

Element	Description	Educational Application
Concise	Clear, focused prompts without unnecessary complexity	Students practice articulating learning goals precisely
Logical	Well-structured prompts with logical flow	Develops organizational and sequencing skills
Explicit	Specific instructions with clear expectations	Requires students to define success criteria
Adaptive	Iterative refinement based on outputs	Builds metacognitive monitoring and adjustment
Restrictive	Appropriate constraints and boundaries	Teaches scope management and critical evaluation

Note: Adapted from Lo [22]. *The Journal of Academic Librarianship*.

- **Concise:** The elimination of linguistic ambiguity, rhetorical flourish, and redundancy to ensure the AI focuses on the core task. This teaches students the value of economy in language, forcing them to distill their intent to its essence. For example, transforming a vague request like “Write something about the Civil War” into “Summarize the three primary economic causes of the American Civil War in under 200 words.”
- **Logical:** The arrangement of instructions in a coherent, step-by-step sequence that the model can process linearly. This mirrors the logic of programming or mathematical proof construction, where the order of operations determines the validity of the result. Students learn to structure their inquiries as a logical chain of thought rather than a disordered stream of consciousness.
- **Explicit:** The clear, granular definition of the desired output format, tone, complexity level, and perspective. This forces the student to demonstrate “genre awareness”—understanding that a scientific abstract requires a fundamentally different tone and structure than a persuasive op-ed. By explicitly defining these parameters for the AI, the student reinforces their own understanding of these genre conventions.
- **Adaptive:** The capacity to refine the prompt based on the initial output, treating the interaction as an iterative dialogue rather than a single transaction. This introduces a metacognitive loop: the student must evaluate the AI’s response against their internal standard of quality, identify the discrepancy, and formulate a corrective instruction. This “meta-prompting” converts the student from a passive consumer to an active editor.
- **Restrictive:** The setting of strict boundaries and negative constraints (e.g., “Do not use sources from before 2020,” “Avoid passive voice,” or “Limit the response to three paragraphs”). This teaches students to preemptively identify potential pitfalls, hallucinations, or irrelevancies, effectively “fencing in” the probability space of the model to ensure accuracy and relevance.

By iterating through these parameters, students learn that the quality of the answer is inextricably contingent upon the precision of the question. This reinforces a fundamental educational truth: that inquiry is a skill as vital as retrieval.

5.2. Structured Scaffolding Approaches

Effective implementation of AI tools requires a trajectory of “fading scaffolding,” a concept rooted in Vygotsky’s Zone of Proximal Development. In the context of AI, this implies a gradual transfer of “prompting responsibility” from the instructor to the learner. In the initial phases, educators should provide “system prompts” or “model inter-

actions” that strictly constrain the AI’s behavior to ensure productive engagement and prevent cognitive bypassing. As students demonstrate competency, ethical judgment, and subject mastery, these constraints are progressively loosened, allowing for greater autonomy and open-ended exploration.

Mollick and Mollick [9] delineate specific pedagogical personas that structure this interaction to achieve distinct cognitive goals. These personas transform the AI from a passive encyclopedia into an active sparring partner, preventing passive consumption and enforcing active cognitive retrieval.

- **The Debate Opponent:** In this mode, the AI is instructed to challenge the student’s thesis with counter-arguments and logical fallacies. This forces the student to defend their reasoning, cite evidence to support their claims, and refine their argumentation. It effectively inoculates students against “echo chambers” by ensuring their ideas are subjected to rigorous, if simulated, scrutiny.
- **The Socratic Tutor:** This persona reverses the traditional flow of information; rather than providing answers, the AI is instructed to probe the student’s understanding through recursive questioning (e.g., “Why do you think that step is correct?” or “How does this connect to the previous concept?”). This scaffolding ensures that the student generates the solution themselves, deepening neural consolidation and preventing the illusion of competence that comes from simply reading an answer.
- **The Devil’s Advocate:** Similar to the debate opponent but focused on identifying weaknesses in logic or evidence, this persona helps students stress-test their ideas before final submission. The AI acts as a “red team,” relentlessly exposing gaps in the student’s reasoning, compelling them to fortify their arguments and anticipate objections.
- **The Simulator:** The AI adopts a specific historical or fictional persona (e.g., “You are a factory worker in 1890s Manchester”), allowing students to practice empathy and perspective-taking within a controlled environment. This transforms abstract historical or literary knowledge into a situated, experiential encounter, requiring students to navigate complex social dynamics and historical constraints.
- **The Mentor/Coach:** Unlike the Socratic Tutor which focuses on content mastery, the Mentor persona focuses on the process of work. It provides feedback on drafts, code snippets, or project plans without doing the work itself. For instance, it might say, “Your introduction is strong, but your transition to the second paragraph is abrupt. Consider how you might link the concept of X to Y.” This supports the student’s executive function and revision skills.

By utilizing these structured roles, educators can ensure that the human-AI interaction remains pedagogically sound, focused on specific learning objectives, and aligned with the principles of active learning.

6. Ethical Considerations and Academic Integrity

Meaningful human-AI collaboration requires a comprehensive re-evaluation of ethical dimensions, extending well beyond the binary of plagiarism. We must address the complex intersection of academic integrity, data sovereignty, and algorithmic bias.

6.1. Academic Integrity as Process Transparency

The rapid advancement of text-generation capabilities has rendered the traditional “policing” model of academic integrity—reliant on detection software—functionally obsolete. Technical countermeasures (AI detectors) suffer from high false-positive rates and are easily circumvented by minor syntactic alterations. Rather than engaging in a futile technological arms race, educators should cultivate a culture of “process transparency.” In the AI era, integrity is defined not by the absence of tools, but by the honest disclosure of their use. Perkins [23] argues for the implementation of “academic transparency statements,” where students are required to formally cite the specific AI prompts utilized and delineate which portions of the work were machine-generated versus human-authored. This practice shifts the evaluative focus from the final product to the “cognitive audit trail” of its creation, ensuring that the student remains the intellectual architect even when leveraging algorithmic assistance.

6.2. Data Privacy and Digital Sovereignty

A critical and often overlooked ethical dimension is the commodification of student data. When students interact with commercial Large Language Models (LLMs), they are frequently, and often unwittingly, feeding proprietary

ideas and personal data into a corporate model's training set. This transforms the learner from a beneficiary of technology into an unpaid data laborer. Educational institutions must therefore establish rigorous guidelines regarding "data sovereignty." This involves the implementation of "walled garden" environments where data is not retained for model training, or the strict enforcement of "PII hygiene" protocols. "Sanitizing" prompts—the systematic removal of names, institutional affiliations, locations, and personal identifiers—must be taught as a foundational digital literacy skill. Students must understand the permanence of their digital footprint within these black-box systems and the potential for their inputs to be regurgitated in other contexts.

6.3. Algorithmic Bias and Epistemic Homogeneity

Furthermore, the uncritical integration of AI poses a risk of "epistemic monoculture." LLMs, predominantly trained on Western, English-centric datasets, inherently encode specific cultural, linguistic, and ideological biases. Reliance on these models can lead to a flattening of cultural nuance and the marginalization of non-dominant dialects or worldviews [24]. There is a profound risk that the "standard English" enforced by these models will be perceived as the only valid form of expression, eroding linguistic diversity. Educators must therefore encourage students to act as "ethical auditors," critiquing the worldview of the AI. Assignments should explicitly invite students to identify where the AI's output reflects a Western-centric bias, stereotypes, or historical omissions, thereby transforming the tool's limitations into a subject of critical inquiry rather than a source of authoritative truth.

7. Implications and Recommendations

7.1. For Teacher Professional Development

It is imperative that professional development (PD) frameworks evolve significantly beyond the rudimentary instruction of tool mechanics—i.e., "how to use the tool"—to address the far more complex domain of "how to teach with the tool." This necessitates a paradigmatic shift from technical training to pedagogical reasoning, where the focus is not on the software itself but on its strategic deployment within specific instructional contexts. Teachers must be afforded dedicated time to experimentally navigate the "jagged frontier" of artificial intelligence—a concept elucidating that AI capabilities are not uniform; the technology may demonstrate superhuman proficiency in certain complex tasks while failing catastrophically in seemingly trivial ones. Understanding this frontier is critical; without this experiential knowledge, educators risk deploying AI in scenarios where it is inherently unreliable or, conversely, failing to leverage it where it could offer substantial support. Furthermore, PD must cultivate the skill of "pedagogical auditing," wherein teachers are trained to critically evaluate AI-generated outputs for subtle biases, factual hallucinations, and pedagogical appropriateness before these outputs are introduced to students. This moves the teacher's role from a passive adopter of technology to an active evaluator and curator of algorithmic content, ensuring that the integration of AI aligns with rigorous educational standards. Recent evidence further indicates that effective implementation depends on faculty self-efficacy, institutional support, and sustained professional development [25,26].

7.2. For Curriculum Design

The advent of generative AI necessitates a fundamental restructuring of curricula to explicitly incorporate AI literacy as a core competency rather than a peripheral elective. This involves a decisive shift in assessment and instructional focus from "answer production"—a domain in which artificial intelligence now excels—to "question formulation" and "output evaluation." In an era where answers are commoditized, the intellectual value migrates to the ability to ask the right questions and to discern the validity of the responses received. Consequently, curricula must prioritize "epistemic humility" and verification skills, compelling students to treat AI outputs not as authoritative truths but as provisional drafts requiring rigorous interrogation. This also implies the integration of "adversarial testing" into the syllabus, where students are tasked with deliberately probing AI systems to identify limitations, biases, and errors. Such exercises transform the relationship with the machine from one of dependence to one of critical mastery. Moreover, this literacy must be interdisciplinary; ethical considerations regarding data privacy, algorithmic bias, and the socio-economic impact of automation should not be sequestered within computer science departments but must be woven into the fabric of humanities, social sciences, and STEM education alike, preparing students for a future where human-AI interaction is ubiquitous. This broader curricular rethinking also aligns

with scholarship that uses ChatGPT as a case study for examining educational chatbot adoption and ethics [27]. Recent reviews also underline that AI literacy, cross-sector implementation, and explicit ethical guidance must be developed systematically rather than assumed [28–30].

8. Case Studies and Practical Applications

Because this manuscript is conceptual rather than a report of a completed intervention study, the following cases are presented as illustrative design scenarios synthesized from recent literature and classroom practice, not as direct empirical proof of the framework. Their purpose is to show how COLLABORATE principles could be operationalized in different contexts and to make future evaluation designs more concrete.

8.1. Case Study 1: The AI as Socratic Tutor in STEM

Application of Principles: Orchestrated Interactions, Learner Agency, Orchestrated Feedback.

In a high school physics curriculum centered on vector mechanics, the instructor instituted a departure from traditional problem sets, replacing them with a mandatory AI-mediated dialogue. This intervention was designed to counteract the “answer-getting” mentality often prevalent in STEM subjects. Utilizing a rigorously engineered system prompt—architecture designed to simulate a “Socratic Tutor” [9]—the artificial intelligence was strictly constrained in its operational parameters.

- **The Constraint Architecture:** The system prompt explicitly forbade the provision of direct solutions or mathematical derivations. Instead, the AI was programmed to respond to student queries solely with guiding questions, analogies, or requests for the student to articulate their own reasoning steps. This constraint forced the interaction into a zone of “productive struggle,” preventing the cognitive offloading of the calculation process.
- **The Outcome and Student Experience:** Initial student reaction was characterized by significant friction; students accustomed to using AI as an answer-retrieval engine attempted to “game” the system, only to be met with persistent refusal and redirective questioning. However, as the semester progressed, this friction transmuted into deeper engagement. Students who accepted the dialogic premise reported a marked increase in conceptual confidence, as they were compelled to verbally derive their logic.
- **Teacher Role Transformation:** The instructor’s role shifted from that of a content deliverer to a “pedagogical data analyst.” By reviewing the anonymized chat logs (the “data exhaust” of the learning process), the teacher could identify widespread misconceptions that the AI failed to adequately address. This allowed for targeted, “Orchestrated Feedback” during in-person lectures, addressing collective blind spots revealed by the algorithmic interaction.

8.2. Case Study 2: Recursive Writing with “Track Changes”

Application of Principles: Assessment Redesign, Reflective Practice, Literacy Development.

In a tertiary-level academic writing course, the assessment protocol was fundamentally redesigned to prioritize the “Process Portfolio” over the final static submission. This intervention sought to leverage AI as a critical editor rather than a content generator, thereby reinforcing the principle of “Reflective Practice.”

- **Phase 1: Unassisted Drafting:** Students were required to produce an initial essay draft within a controlled, technology-free environment (“Balanced Integration”). This ensured that the foundational argumentation and voice were authentically human.
- **Phase 2: Algorithmic Critique:** Students then submitted their drafts to an AI model, prompting it with specific, rubric-aligned criteria (e.g., “Analyze the logical coherence of my counter-arguments” or “Identify areas where evidence is insufficient”). This phase utilized the AI as a high-availability feedback mechanism.
- **Phase 3: The Editorial Defense:** Crucially, students were not required to accept the AI’s suggestions. The final submission included the original draft, the AI’s critique, the revised manuscript, and—most importantly—a “defense statement.” In this statement, students had to explicitly justify why they accepted certain AI edits and why they rejected others.
- **Result:** This structure validated the premise that “refusal is a skill.” Students demonstrated critical thinking not by utilizing the tool, but by overruling it. The assessment measured the student’s ability to maintain intellectual sovereignty in the face of algorithmic persuasion, shifting academic integrity from a policing model to

a transparency model.

8.3. Case Study 3: Historical Simulation and Perspective Taking

Application of Principles: Contextualized Learning, Authentic Tasks.

In an advanced history curriculum, the principle of “Contextualized Learning” was operationalized through an AI-driven simulation. Students were tasked with interrogating a “persona” of a historical figure, generated by the AI, to explore complex geopolitical decisions.

- **The Simulation Design:** Rather than writing a standard biographical essay, students engaged in a real-time debate with an AI instructed to adopt the worldview, linguistic style, and knowledge limitations of a specific historical actor (e.g., a Cold War diplomat). The prompt engineering required the AI to avoid modern hindsight bias.
- **The Assessment of Interrogation:** Students were assessed not on the AI’s output, but on the quality of their interrogation strategies. Did they ask questions that revealed the persona’s underlying biases? Did they successfully identify historical anachronisms or “hallucinations” in the AI’s responses?
- **The Outcome:** This task transformed the AI from a source of information into an object of scrutiny. It required students to possess deep subject matter knowledge to effectively challenge the simulation, thereby reinforcing the necessity of human expertise. The activity underscored that in a world of synthetic media, the ultimate skill is the ability to verify and contextualize automated narratives.

These scenarios are analytically useful, but they do not substitute for controlled implementation research, expert validation, or longitudinal evidence. Future studies should therefore test whether process transparency, AI-free phases, structured feedback roles, and discipline-specific orchestration actually improve learning outcomes, reduce inappropriate AI dependence, and preserve foundational skills across different cultural and institutional settings. Recent reviews of assessment transformation and pedagogical implementation likewise reinforce the need for more rigorous, longitudinal, and context-sensitive evidence [31,32].

9. Limitations and Future Research

Several limitations must be acknowledged explicitly. First, COLLABORATE is a conceptual synthesis and has not yet been empirically validated as an integrated intervention. Second, much of the evidence based on classroom AI remains concentrated in higher education and short-duration studies, constraining claims about long-term effectiveness. Third, institutional resources, language background, policy environments, and cultural expectations may influence how the framework is interpreted and enacted. Future work should therefore include quasi-experimental classroom studies, cross-cultural comparisons, expert review of principle coherence, and mixed-method analysis of how different collaboration models affect metacognition, revision quality, and academic integrity. Recent meta-reviews and higher education reviews likewise show that the field remains fast-moving, heterogeneous, and still short on strong longitudinal evidence across contexts [33–35].

10. Conclusions

The integration of generative artificial intelligence into education has made an old assumption untenable: a polished product can no longer be treated as sufficient evidence of learning. The challenge is therefore not simply whether AI should be allowed, but how teachers can redesign pedagogy so that human judgment, disciplinary reasoning, and ethical responsibility remain central. The revised version of this manuscript argues that this redesign requires more than policy statements or prompt tips; it requires a coherent pedagogical framework that links context, orchestration, agency, feedback, and assessment.

COLLABORATE is offered as that framework in conceptual form. Its value lies not in claiming final empirical proof, but in providing a traceable design logic for moving from automation toward augmentation under real classroom constraints. By making process transparency, reflective use, authentic tasks, and teacher empowerment central, the framework responds directly to the core risks identified in the introduction: invisible outsourcing, cognitive offloading, and assessment misalignment. The next step is empirical: to test which combinations of these principles most effectively sustain foundational skills while enabling responsible and meaningful human-AI collaboration.

Funding

This research received no external funding.

Institutional Review Board Statement

Not applicable. This article is a conceptual review and did not involve human or animal participants.

Informed Consent Statement

Not applicable.

Data Availability Statement

No new data were created or analyzed in this study. Data sharing is not applicable to this article.

Acknowledgments

The author thanks the editors and reviewers for their constructive feedback on earlier versions of this manuscript.

Conflicts of Interest

The author declares no conflict of interest.

AI Use Statement

During the preparation of this work, the author used ChatGPT for language refinement and editorial revision. The author subsequently reviewed and edited the content as necessary and takes full responsibility for the final content of the published article.

References

1. Dwivedi, Y.K.; Kshetri, N.; Hughes, L.; et al. So what if ChatGPT wrote it? Multidisciplinary perspectives on opportunities, challenges and implications of generative conversational AI for research, practice and policy. *Int. J. Inf. Manag.* **2023**, *71*, 102642.
2. Kasneci, E.; Sessler, K.; Küchemann, S.; et al. ChatGPT for good? On opportunities and challenges of large language models for education. *Learn. Individ. Differ.* **2023**, *103*, 102274.
3. UNESCO. *UNESCO Survey: Less than 10% of Schools and Universities Have Formal Guidance on AI*; UNESCO: Paris, France, 2023.
4. Freeman, J. *Student Generative AI Survey 2025*; Higher Education Policy Institute: Oxford, UK, 2025.
5. Trust, T.; Whalen, J.; Mouza, C. Editorial: ChatGPT: Challenges, Opportunities, and Implications for Teacher Education. *Contemp. Issues Technol. Teach. Educ.* **2023**, *23*, 1–23.
6. Lodge, J.M.; Howard, S.K.; Bearman, M.; et al. Assessment reform for the age of artificial intelligence. *Tert. Educ. Manag.* **2023**, 1–12. Available online: <https://www.teqsa.gov.au/sites/default/files/2023-09/assessment-reform-age-artificial-intelligence-discussion-paper.pdf>
7. Scarfe, P.; Watcham, K.; Clarke, A.; et al. A real-world test of artificial intelligence infiltration of a university examinations system: A ‘Turing Test’ case study. *PLoS ONE* **2024**, *19*, e0305354. [CrossRef]
8. Wilson, H.J.; Daugherty, P.R. Collaborative intelligence: Humans and AI are joining forces. *Harv. Bus. Rev.* **2018**, *96*, 114–123.
9. Mollick, E.; Mollick, L. Assigning AI: Seven Approaches for Students, with Prompts. *SSRN Electron. J.* **2023**, 1–48. Available online: https://papers.ssrn.com/sol3/papers.cfm?abstract_id=4475995
10. Dellermann, D.; Ebel, P.; Söllner, M.; et al. Hybrid Intelligence. *Bus. Inf. Syst. Eng.* **2019**, *61*, 637–643.
11. Holstein, K.; McLaren, B.M.; Alevan, V. Co-Designing a Real-Time Classroom Orchestration Tool to Support Teacher–AI Complementarity. *J. Learn. Anal.* **2019**, *6*, 27–52.
12. Sharples, M.; Pérez y Pérez, R. *Story Machines: How Computers Have Become Creative Writers*; Routledge: London, UK, 2022.
13. Prather, J.; Reeves, B.N.; Denny, P.; et al. “It’s Weird That It Knows What I Want”: Usability and Interactions with Copilot for Novice Programmers. *ACM Trans. Comput.-Hum. Interact.* **2023**, *31*, 4.

14. Gero, K.I.; Ashktorab, Z.; Dugan, C.; et al. Mental Models of AI Agents in a Cooperative Word Guessing Game. In Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems, Honolulu, HI, USA, 25–30 April 2020; pp. 610–623.
15. Attewell, S. *Student Perceptions of AI 2025*; Jisc: Bristol, UK, 2025.
16. McDonald, N.; Johri, A.; Ali, A.; et al. Generative artificial intelligence in higher education: Evidence from an analysis of institutional policies and guidelines. *Comput. Hum. Behav. Artif. Hum.* **2025**, *3*, 100121. [[CrossRef](#)]
17. Bond, M.; Khosravi, H.; De Laat, M.; et al. A meta-systematic review of artificial intelligence in higher education: A call for increased ethics, collaboration, and rigour. *Int. J. Educ. Technol. High. Educ.* **2024**, *21*, 4. [[CrossRef](#)]
18. Banihashem, S.K.; Noroozi, O.; Khosravi, H.; et al. Pedagogical framework for hybrid intelligent feedback. *Innov. Educ. Teach. Int.* **2025**, *63*, 554–570. [[CrossRef](#)]
19. Banihashem, S.K.; Bond, M.; Khosravi, H.; et al. A systematic mapping review at the intersection of artificial intelligence and self-regulated learning. *Int. J. Educ. Technol. High. Educ.* **2025**, *22*, 50. [[CrossRef](#)]
20. Chan, C.K.Y. A comprehensive AI policy education framework for university teaching and learning. *Int. J. Educ. Technol. High. Educ.* **2023**, *20*, 38. [[CrossRef](#)]
21. Topping, K.J.; Gehringer, E.; Khosravi, H.; et al. Enhancing peer assessment with artificial intelligence. *Int. J. Educ. Technol. High. Educ.* **2025**, *22*, 3. [[CrossRef](#)]
22. Lo, L.S. The CLEAR path: A framework for enhancing information literacy through prompt engineering. *J. Acad. Librariansh.* **2023**, *49*, 102720.
23. Perkins, M. Academic Integrity in the Age of AI. *Nat. Mach. Intell.* **2023**, *5*, 2–4.
24. Bender, E.M.; Gebru, T.; McMillan-Major, A.; et al. On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? In Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency, Online, 3–10 March 2021; pp. 610–623.
25. Mah, D.-K.; Groß, N. Artificial intelligence in higher education: Exploring faculty use, self-efficacy, distinct profiles, and professional development needs. *Int. J. Educ. Technol. High. Educ.* **2024**, *21*, 58. [[CrossRef](#)]
26. Ifenthaler, D.; Majumdar, R.; Gorissen, P.; et al. Artificial Intelligence in Education: Implications for Policy-makers, Researchers, and Practitioners. *Technol. Knowl. Learn.* **2024**, *29*, 1693–1710. [[CrossRef](#)]
27. Tlili, A.; Shehata, B.; Adarkwah, M.A.; et al. What if the devil is my guardian angel: ChatGPT as a case study of using chatbots in education. *Smart Learn. Environ.* **2023**, *10*, 15. [[CrossRef](#)]
28. Yim, I.H.Y.; Su, J. Artificial intelligence (AI) learning tools in K-12 education: A scoping review. *J. Comput. Educ.* **2025**, *12*, 93–131. [[CrossRef](#)]
29. Biagini, G. Towards an AI-Literate Future: A Systematic Literature Review Exploring Education, Ethics, and Applications. *Int. J. Artif. Intell. Educ.* **2025**, *35*, 2616–2666. [[CrossRef](#)]
30. Memarian, B.; Doleck, T. Teaching and learning artificial intelligence: Insights from the literature. *Educ. Inf. Technol.* **2024**, *29*, 21523–21546. [[CrossRef](#)]
31. Xia, Q.; Weng, X.; Ouyang, F.; et al. A scoping review on how generative artificial intelligence transforms assessment in higher education. *Int. J. Educ. Technol. High. Educ.* **2024**, *21*, 40. [[CrossRef](#)]
32. Qian, Y. Pedagogical Applications of Generative AI in Higher Education: A Systematic Review of the Field. *TechTrends* **2025**, *69*, 1105–1120. [[CrossRef](#)]
33. Fu, Y.; Weng, Z.; Wang, J. Examining AI Use in Educational Contexts: A Scoping Meta-Review and Bibliometric Analysis. *Int. J. Artif. Intell. Educ.* **2025**, *35*, 1388–1444. [[CrossRef](#)]
34. Jensen, L.X.; Buhl, A.; Sharma, A.; et al. Generative AI and higher education: A review of claims from the first months of ChatGPT. *High. Educ.* **2025**, *89*, 1145–1161. [[CrossRef](#)]
35. McGrath, C.; Farazouli, A.; Cerratto-Pargman, T. Generative AI chatbots in higher education: A review of an emerging research area. *High. Educ.* **2025**, *89*, 1533–1549. [[CrossRef](#)]



Copyright © 2025 by the author(s). Published by UK Scientific Publishing Limited. This is an open access article under the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

Publisher’s Note: The views, opinions, and information presented in all publications are the sole responsibility of the respective authors and contributors, and do not necessarily reflect the views of UK Scientific Publishing Limited and/or its editors. UK Scientific Publishing Limited and/or its editors hereby disclaim any liability for any harm or damage to individuals or property arising from the implementation of ideas, methods, instructions, or products mentioned in the content.