

CAT-AI: Supporting Teacher Workflows with AI-Assisted Exercise Creation

*Original*

CAT-AI: Supporting Teacher Workflows with AI-Assisted Exercise Creation / Calò, Tommaso; Lorenzo, Cuccu; De Russis, Luigi. - ELETTRONICO. - (2026), pp. 1-6. ( CHI '26: CHI Conference on Human Factors in Computing Systems Barcelona (ESP) 13–17 April 2026) [10.1145/3772363.3799367].

*Availability:*

This version is available at: 11583/3008331 since: 2026-04-24T12:08:06Z

*Publisher:*

ACM

*Published*

DOI:10.1145/3772363.3799367

*Terms of use:*

This article is made available under terms and conditions as specified in the corresponding bibliographic description in the repository

*Publisher copyright*

(Article begins on next page)

# CAT-AI: Supporting Teacher Workflows with AI-Assisted Exercise Creation

Tommaso Calò  
Dipartimento di Automatica e  
Informatica  
Politecnico Di Torino  
Torino, Torino, Italy  
tommaso.calo@polito.it

Lorenzo Cuccu  
Politecnico di Torino  
Torino, Italy  
s330260@studenti.polito.it

Luigi De Russis  
Dipartimento di Automatica e  
Informatica  
Politecnico di Torino  
Torino, Italy  
luigi.derussis@polito.it

## Abstract

Creating classroom exercises consumes substantial teacher time, requiring educators to balance pedagogical soundness with student engagement while adapting materials for diverse learners. Through formative interviews with ten K-12 teachers, we observed that educators increasingly turn to AI tools for exercise creation, yet struggle with prompt design, fragmented workflows across multiple applications, and significant verification overhead. From these interviews we synthesized a three-phase workflow model describing how teachers conceptualize, transform, and finalize educational content. Building on these insights, we present CAT-AI, an AI-assisted authoring system that embeds structured parameter specification, unified WYSIWYG editing, and transparent confidence indicators to support teachers throughout exercise creation without requiring prompt engineering expertise, disrupting workflow, or extensive manual verification.

## CCS Concepts

• **Human-centered computing** → **Human computer interaction (HCI)**; **Interactive systems and tools**; • **Applied computing** → **Education**.

## Keywords

Education and Learning Technologies, Generative AI, Teacher-AI Collaboration, Human-Centered AI Design

### ACM Reference Format:

Tommaso Calò, Lorenzo Cuccu, and Luigi De Russis. 2026. CAT-AI: Supporting Teacher Workflows with AI-Assisted Exercise Creation. In *Extended Abstracts of the 2026 CHI Conference on Human Factors in Computing Systems (CHI EA '26)*, April 13–17, 2026, Barcelona, Spain. ACM, New York, NY, USA, 6 pages. <https://doi.org/10.1145/3772363.3799367>

## 1 Introduction

Educational content creation consumes substantial teacher time, with educators spending hours adapting materials for diverse learners while balancing pedagogical objectives, difficulty calibration, and student engagement [14, 20]. Studies reveal that teachers employ transformation-based approaches to exercise creation, using techniques like addition of supplementary materials, deletion of

inappropriate content, and modification of difficulty levels [15, 17]. Traditionally, digital efforts to support teachers have focused on Intelligent Tutoring Systems (ITS) and related authoring environments, which enabled educators to design and train digital tutors without programming expertise [3, 4]. However, their primary focus was on the digitization of tutoring practices rather than on pedagogical content creation.

More recently, the emergence of AI-powered tools has shifted attention directly to educational content creation, supporting teachers in generating and adapting classroom exercises, worksheets, and assessments. Early neural network approaches demonstrated technical feasibility for generating questions from textbooks [10, 18], and LLMs promised greater flexibility through prompt engineering [11], with validated capabilities now integrated into commercial platforms [6, 13].

Han et al. [8] provide empirical evidence that while teachers appreciate AI's ability to accelerate material preparation, they face a significant verification burden: most Generative AI systems provide no indication of output reliability, forcing teachers to manually check every element of the generated content [5, 7]. This lack of transparency undermines adoption, as the effort spent verifying outputs can offset the time saved through generation. Separately, misalignment with pedagogical goals often derives from prompt engineering difficulties, as teachers must translate nuanced instructional intentions into textual instructions [16, 19]. Template-based systems like REDEEM helped non-programmer teachers create interactive content through structured interfaces [1, 2]. However, these systems assume teachers start from scratch, offering limited possibilities to reuse or adapt existing resources. On the contrary, field studies show teachers using AI-generated content as starting points, extensively modifying outputs for their specific contexts [8, 9]. This transformation-based approach represents teachers' natural workflow but cannot be easily expressed through current interfaces, which fragment generation, editing, and formatting into disconnected stages [12].

In this work, we examine how teachers engage with AI throughout their authoring workflow. We conducted formative interviews with 10 K-12 educators using semi-structured protocols and think-aloud observations. Teachers demonstrated their typical exercise creation workflows, revealing recurring challenges that cluster around four areas: difficulty translating pedagogical requirements into effective prompts; use of multiple disconnected tools creating substantial context-switching overhead; effort spent verifying AI-generated content without transparency about reliability; and



This work is licensed under a Creative Commons Attribution 4.0 International License. *CHI EA '26, Barcelona, Spain*

© 2026 Copyright held by the owner/author(s).  
ACM ISBN 979-8-4007-2281-3/26/04  
<https://doi.org/10.1145/3772363.3799367>

inability to easily specify what to preserve versus modify when adapting existing materials.

These observations led us to formalize a three-phase workflow model. During conceptualization, teachers specify pedagogical parameters through structured fields rather than open-ended prompts, eliminating the need to translate teaching expertise into prompt engineering. During transformation, a unified What You See Is What You Get (WYSIWYG) editor enables both direct text manipulation and targeted AI-assisted regeneration, with confidence indicators flagging potentially problematic content. During finalization, the system handles export formatting automatically, producing classroom-ready materials without requiring external tools. We present CAT-AI, a web-based system that operationalizes teachers' workflow through integrated support for each phase.

CAT-AI is available as open-source software with a live demonstration.<sup>1</sup>

## 2 Formative Study

To ground our design in observed pedagogical practices, we conducted semi-structured remote interviews with ten K-12 teachers, aimed at understanding their current workflows for creating classroom exercises and identifying friction points in AI-assisted content creation. We recruited teachers from primary and secondary schools covering grades 1 through 8, with teaching experience ranging from 3 to 40 years and varying levels of AI tool proficiency, balanced between primary (n=5) and middle school (n=5). Interviews lasted approximately 20 minutes each and were recorded with participant consent. The interview protocol covered five areas: (1) recent teaching activities and exercise types created; (2) tools, processes, and thinking behind exercise creation; (3) perceived advantages and disadvantages of digital versus paper-based materials; (4) strategies for adapting exercises to heterogeneous student abilities; and (5) collection of example artifacts such as exercise files and photographs of recently created materials. Interview transcripts and collected artifacts were analyzed qualitatively to inform both the workflow model (Figure 1) and the design goals described in Section 3.

### 2.1 Findings

Across all participants, we identified a common multi-tool workflow (Figure 1) characterized by four recurring stages: planning exercises around learning objectives, calibrating difficulty to student abilities, generating and verifying content through AI tools, and adapting existing materials rather than creating from scratch. Each stage surfaced specific friction points that informed our design.

**Planning.** Teachers begin by identifying learning objectives from their curriculum and determining target difficulty levels based on their knowledge of student abilities. All interviewed teachers reported that learning objectives drive exercise design. Clarity in instructions is critical since ambiguous wording undermines the purpose of an exercise. Coherence emerged as equally important across three dimensions: terminological coherence requires consistent use of the vocabulary introduced in class, structural coherence means following familiar exercise patterns, and visual coherence demands consistent formatting conventions.

**Difficulty calibration.** Adjusting exercise complexity to match student ability was universally cited as critical but context-dependent. Primary teachers emphasized that all students should be able to complete exercises successfully, while middle school teachers focused on testing higher-order thinking rather than knowledge recall. Teachers stressed that appropriate difficulty can only be determined through direct knowledge of individual students.

**Content generation and verification.** The generation phase, in which teachers create exercises, worksheets, and assessment materials, exposes challenges in translating pedagogical intent into effective prompts. Teachers navigate to ChatGPT or similar AI tools and construct prompts through iterative refinement, yet the gap between pedagogical expertise and prompt engineering skills creates repeated cycles of adjustment. As one teacher explained, "I have to be very specific about what I want, but even then the system frequently ignores these constraints" (T6). Once content is generated, teachers manually transfer it into word processors, introducing workflow discontinuity. Teachers then extensively modify materials by simplifying instructions, adjusting vocabulary, correcting errors, and adding visual elements. One participant noted: "I spend nearly as much time checking and fixing generated content as I would creating from scratch" (T8).

**Adaptation workflow.** Throughout the study, teachers rarely create exercises from scratch but instead take existing materials and selectively modify them through vocabulary adjustment, content substitution, and difficulty modification. As T6 described, "I find an exercise that's close to what I need, then change the parts that don't fit." Current AI tools do not support this transformation-based workflow: they require teachers to articulate requirements through prompt engineering, they operate through holistic regeneration that produces entirely new content, they provide no indication of output reliability forcing extensive manual verification, and they exist as disconnected applications that fragment the creation process across generation, editing, and formatting tools.

## 3 Design Goals

While generative models offer powerful content creation capabilities, the formative study revealed evident frictions in how teachers currently use AI for exercise creation. The cognitive overhead of prompt engineering, workflow fragmentation, and verification burden often negate potential efficiency gains. These findings informed three design goals for CAT-AI.

**Goal 1: Minimize Prompt Engineering Burden.** Teachers possess deep pedagogical expertise, yet current tools require them to develop prompt engineering skills to leverage AI assistance. Rather than asking teachers to translate requirements into natural language prompts, the system should directly capture dimensions teachers naturally consider: learning objectives, prerequisites, school level, and difficulty calibration. The dual-mode generation approach, allowing either file-based reference or manual specification, emerged from interviews where teachers rarely create from scratch but work from existing materials.

<sup>1</sup>Platform demo: <https://catai.paperbackwriters.club/>; Code: <https://repo.paperbackwriters.club/code/catai>.

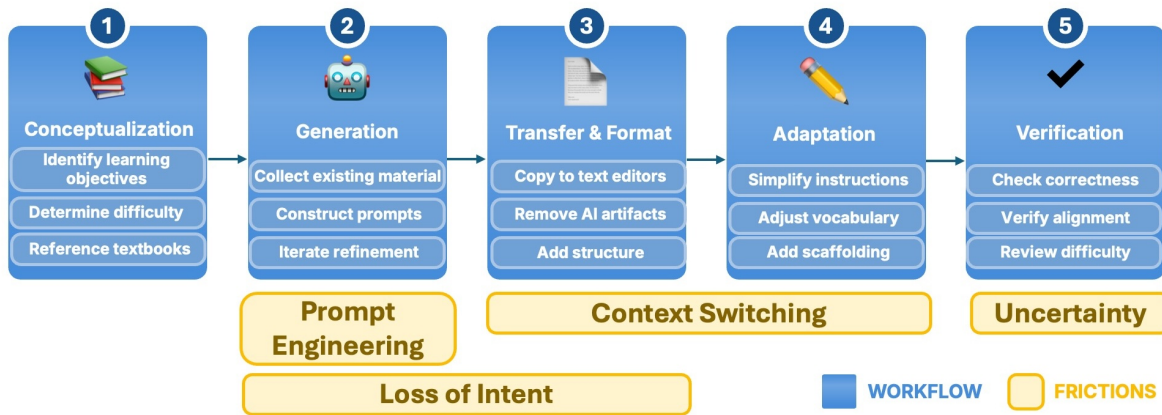


Figure 1: Typical teacher workflow for AI-assisted exercise creation, showing fragmentation across multiple tools.

### Goal 2: Integrate Generation, Customization, and Export.

The observed workflow fragmentation creates significant cognitive overhead as teachers move between ChatGPT for generation, Word for formatting, and various verification steps. A unified environment should support WYSIWYG editing where teachers see content in its final printed form, provide granular control over AI-assisted modifications allowing selective regeneration while preserving context, and handle export formatting automatically to produce classroom-ready materials without manual reformatting.

**Goal 3: Support Verification Through Transparency.** Teachers reported spending considerable effort checking AI-generated content because existing tools provide no indication of output reliability. To reduce this verification burden, the system should include confidence indicators that flag potentially problematic content, directing teacher attention to elements most likely to require scrutiny rather than forcing exhaustive review. Multiple regeneration modalities, from complete exercise regeneration to element-level adjustments, support teachers' transformation-based workflow by distinguishing what should change from what must remain constant.

## 4 CAT-AI System

We implemented CAT-AI as a full-stack web application that translates our design goals into a functional system supporting teachers' three-phase workflow (Figure 2). The frontend uses React 18 with Material-UI components, while the backend uses Node.js with Express, integrating with OpenAI's GPT-4o for content generation and verification.

The workflow begins with a **conceptualization phase** where teachers specify pedagogical parameters without requiring prompt construction (Figure 2, top row). From the welcome screen, teachers select between file-based reference or manual entry modes. In file mode, the system accepts PDF or image uploads as shown in the file upload panel, including screenshots or photos of existing exercises, extracts text content server-side using pdf-parse for PDFs or OCR for images, and automatically infers candidate learning objectives, prerequisites, school level, and grade through an initial LLM API call. Teachers review and modify these extracted parameters before

generation proceeds. Manual mode presents the structured input forms visible in the manual input panel, directly requesting pedagogical dimensions such as school level, grade, learning objectives as tagged entries, prerequisites, and optional features like worked examples or student reminders. Both modes converge to a common generation pipeline through the pedagogical settings interface.

Generation settings including difficulty level, number of exercises, vocabulary tone, language, and visual style apply globally across sessions through a persistent settings modal. This separation between exercise-specific parameters, representing what to teach, and stylistic preferences, representing how to present, reduces cognitive load during individual exercise creation. The generation process involves two sequential API calls. The first generates exercise and solution content structured as JSON arrays of text elements, each with content, positioning, dimensions, and styling attributes; the second call analyzes this generated content for logical or mathematical errors, assigning confidence scores and explanatory notes to each element. This dual-pass approach is designed to mitigate the reliability concerns raised by teachers, as the system explicitly flags uncertain content rather than presenting all output with equal authority.

Once content is generated, teachers enter the **transformation phase**. Through a WYSIWYG editor they directly interact with exercise content in its final printed form (Figure 2, middle row). The direct manipulation interface eliminates the abstraction gap between editing and output that characterizes the ChatGPT and Microsoft Word workflows observed in our formative study. As visible in the editor interface, exercise and solution occupy separate A4-sized canvases where all content appears as draggable, resizable text or image elements.

The editor supports three distinct modification modalities. Manual editing allows direct text modification through double-click interaction, drag-and-drop repositioning with alignment guides, element resizing within page boundaries, and styling adjustments including font size, color, and emphasis. The AI-assisted editing panel enables element-level regeneration where teachers select

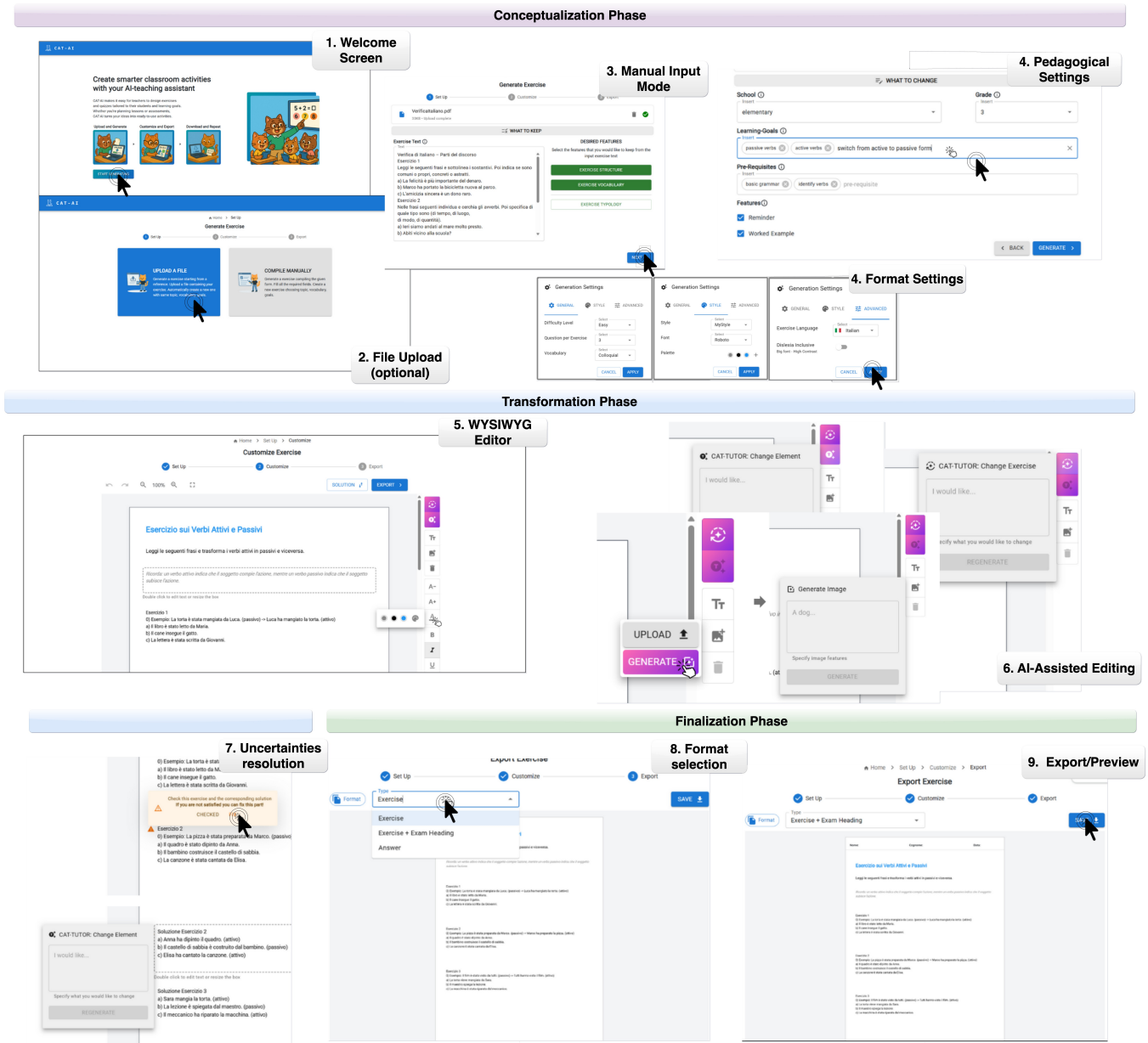


Figure 2: Complete user interaction flow through CAT-AI’s three-phase workflow. The conceptualization phase (top) guides structured parameter specification through welcome screen, optional file upload, manual input, and pedagogical settings. The transformation phase (middle) provides uncertainty resolution, AI-assisted editing, and a WYSIWYG editor interface. The finalization phase (bottom) handles format selection and export with preview.

any text element and request modifications through natural language prompts, with the system preserving all other content while regenerating only the specified element. Exercise-level AI regeneration processes natural language requests to add, modify, or remove content across the entire exercise or solution, supporting the transformation-based workflow where teachers make targeted changes to mostly-satisfactory content. The uncertainty resolution interface, shown in the middle row, presents flagged content

with confidence indicators and one-click “Fix” options that trigger targeted regeneration.

Additional AI-assisted features include solution recalculation that regenerates solutions when exercise content changes to maintain consistency, and image generation through DALL-E 2 with text prompts. Each AI operation maintains context from the current exercise state, ensuring modifications remain coherent with existing content rather than producing disconnected alternatives.

The implementation stores all element data in React Context as structured objects containing position, dimensions, content, styling, confidence scores, and unique identifiers, enabling immediate UI updates and undo/redo functionality.

The **finalization phase** addresses the formatting burden observed in current workflows where teachers spend considerable time reformatting content for distribution (Figure 2, bottom row). Through the format selection interface, teachers preview their complete exercise in standard A4 layout and select from predefined formats: exercise only, exercise with exam header including name, date, and class fields, or solution only. The export panel converts the React component representation to PDF server-side, handling page margins, element positioning, and formatting automatically. Generated PDFs become immediately available for download without additional processing, eliminating the manual reformatting work that fragments teachers' current workflows.

The system architecture separates presentation logic from content generation and file management. The frontend employs React 18<sup>2</sup> for component architecture, Material-UI<sup>3</sup> for interface components, and react-draggable<sup>4</sup> for element positioning. PDF generation uses html2canvas<sup>5</sup> and jsPDF<sup>6</sup>. The backend runs on Node.js<sup>7</sup> with Express<sup>8</sup> for routing, Multer<sup>9</sup> for PDF text extraction, and Tesseract.js<sup>10</sup> for OCR on image uploads. Content generation and verification use OpenAI's GPT-4o<sup>11</sup>, while image generation uses DALL-E 2<sup>12</sup>. Users provide their own API key, which is stored only in the browser session and transmitted over HTTPS without server-side persistence. The dual-pass generation pipeline introduces a latency of approximately 50 seconds per exercise, a trade-off between verification coverage and responsiveness.

## 5 Conclusions

CAT-AI is a prototype system that attempts to align AI assistance with teachers' natural exercise creation practices instead of requiring educators to adapt to general-purpose tools. The system operationalizes a three-phase workflow model derived from formative interviews, but has not yet been validated through comparative evaluation with target users.

Several limitations should be noted. The confidence indicator mechanism relies on a second-pass analysis using the same language model that generated the content. This means the system may systematically miss certain error classes, producing confident-but-wrong outputs, or flag correct content as uncertain. Teachers should interpret confidence indicators as attention guidance rather than guarantees of correctness. Practical deployment in K-12 settings raises additional considerations: the dual-pass generation pipeline introduces latency of approximately 50 seconds per exercise, API costs accumulate with repeated regeneration, and the

system depends on stable internet connectivity and access to a commercial API. Because exercise content is sent to an external API for generation and verification, schools would need to establish data governance policies, particularly regarding what content may be processed through third-party services.

Several directions for future work emerge from this design. Comparative evaluations with teachers using both CAT-AI and their current practices would measure whether the structured approach reduces cognitive load and improves content quality. Longitudinal deployment would reveal how teachers integrate the system into regular practice and whether initial benefits persist over time. Future versions might incorporate formal knowledge representations such as curriculum standards to provide more pedagogically grounded validation. Detecting and correcting for potential biases in AI-generated educational content represents another important direction. Cross-cultural evaluation spanning different countries and school systems would clarify which design principles generalize and which require contextual adaptation.

## References

- [1] Shaaron Ainsworth, Nigel Major, Shirley Grimshaw, Mary Hayes, Jean Underwood, Ben Williams, and David Wood. 2003. REDEEM: Simple Intelligent Tutoring Systems from Usable Tools. In *Authoring Tools for Advanced Technology Learning Environments*, Tom Murray, Stephen B. Blessing, and Shaaron Ainsworth (Eds.). Springer, Dordrecht, 205–232. doi:10.1007/978-94-017-0819-7\_8
- [2] Shaaron E. Ainsworth, Shirley K. Grimshaw, and D. Jean Underwood. 1999. Teachers as designers: Using REDEEM to create ITSs for the classroom. *Computers & Education* 33, 2-3 (1999), 171–188.
- [3] Vincent Alevan, Bruce M. McLaren, Jonathan Sewall, and Kenneth R. Koedinger. 2006. The Cognitive Tutor Authoring Tools (CTAT): Preliminary Evaluation of Efficiency Gains. In *Intelligent Tutoring Systems*. Springer, 61–70. doi:10.1007/11774303\_7
- [4] Vincent Alevan, Bruce M. McLaren, Jonathan Sewall, Martin van Velsen, Octav Popescu, Sandra Demi, Michael Ringenberg, and Kenneth R. Koedinger. 2016. Example-Tracing Tutors: Intelligent Tutor Development for Non-Programmers. *Int. J. Artif. Intell. Educ.* 26, 1 (2016), 224–269. doi:10.1007/s40593-015-0088-2
- [5] Deepak Varuvel Dennison, Bakhtawar Ahtisham, Kavyansh Chourasia, Nirmitt Arora, Rahul Singh, Rene F. Kizilcec, Akshay Nambi, Tanuja Ganu, and Aditya Vashistha. 2025. Teacher-AI Collaboration for Curating and Customizing Lesson Plans in Low-Resource Schools. arXiv:2507.00456 [cs.CY].
- [6] Kristen DiCerbo. 2023. Khan Academy explores the potential for GPT-4 in a limited pilot program. (2023). <https://openai.com/customer-stories/khan-academy>.
- [7] Jaimie Drozdal, Justin D Weisz, Dakuo Wang, Gaurav Dass, Bingsheng Yao, Changruo Zhao, Michael Muller, Lin Ju, and Hui Su. 2020. Trust in AutoML: Exploring Information Needs for Establishing Trust in Automated Machine Learning Systems. In *Proceedings of the 25th International Conference on Intelligent User Interfaces*. 297–307. doi:10.1145/3377325.3377501
- [8] Yueqi Han, Yuling Zhou, Hanqi Cai, et al. 2024. Teachers, Parents, and Students' Perspectives on Integrating Generative AI into Elementary Literacy Education. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*. doi:10.1145/3613904.3642438
- [9] Tae Soo Kim, Yoonjoo Lee, Jamin Shin, Young-Ho Kim, and Juho Kim. 2024. EvalLM: Interactive Evaluation of Large Language Model Prompts on User-Defined Criteria. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. 1–21. doi:10.1145/3613904.3642216
- [10] Ekaterina Kochmar, Dung Do Vu, Robert Belfer, Varun Gupta, Iulian Vlad Serban, and Joelle Pineau. 2021. Automated Data-Driven Generation of Personalized Pedagogical Interventions in Intelligent Tutoring Systems. *International Journal of Artificial Intelligence in Education* 32 (2021), 323–349. doi:10.1007/s40593-021-00267-x
- [11] Jaewook Lee and Sanghoon Kim. 2024. Leveraging ChatGPT for Adaptive Learning through Personalized Prompt-based Instruction: A CS1 Education Case Study. In *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems*. doi:10.1145/3613905.3637148
- [12] Xinyi Lu, Simin Fan, Jessica Houghton, Lu Wang, and Xu Wang. 2023. ReadingQuizMaker: A Human-NLP Collaborative System that Supports Instructors to Design High-Quality Reading Quiz Questions. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–18. doi:10.1145/3544548.3580957

<sup>2</sup><https://react.dev>

<sup>3</sup><https://mui.com>

<sup>4</sup><https://github.com/react-grid-layout/react-draggable>

<sup>5</sup><https://html2canvas.hertzen.com>

<sup>6</sup><https://github.com/parallax/jsPDF>

<sup>7</sup><https://nodejs.org>

<sup>8</sup><https://expressjs.com>

<sup>9</sup><https://github.com/expressjs/multer>

<sup>10</sup><https://tesseract.projectnaptha.com>

<sup>11</sup><https://openai.com/index/hello-gpt-4o/>

<sup>12</sup><https://openai.com/dall-e-2>

- [13] MagicSchool ai. 2023. AI for teachers - lesson planning and more! (2023). <https://www.magicschool.ai/>.
- [14] Jason K McDonald. 2021. The Everydayness of Instructional Design and the Pursuit of Quality in Online Courses. *Online Learning* 25, 4 (2021), 156–173.
- [15] Ni Nyoman Padmadewi and Luh Putu Artini. 2024. Textbook Adaptation Techniques in a Technology-Integrated Environment by an Indonesian EFL Teacher. *TEFLIN Journal* 35, 1 (2024), 89–107.
- [16] Auste Simkute, Viktor Kewenig, Abigail Sellen, Sean Rintel, and Lev Tankelevitch. 2025. The New Calculator? Practices, Norms, and Implications of Generative AI in Higher Education. *arXiv preprint arXiv:2501.08864* (2025).
- [17] Brian Tomlinson. 2022. Theorising Textbook Adaptation in English Language Teaching. *Innovation in Language Learning and Teaching* 16, 3 (2022), 218–231.
- [18] Zichao Wang, Andrew S Lan, Weili Nie, Andrew E Waters, Phillip J Grimaldi, and Richard G Baraniuk. 2018. QG-net: A Data-Driven Question Generation Model for Educational Content. In *Proceedings of the Fifth Annual ACM Conference on Learning at Scale*. 1–10. doi:10.1145/3231644.3231654
- [19] J.D. Zamfirescu-Pereira, Richmond Y. Wong, Bjoern Hartmann, and Qian Yang. 2023. Why Johnny Can't Prompt: How Non-AI Experts Try (and Fail) to Design LLM Prompts. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–21. doi:10.1145/3544548.3581388
- [20] Xiaofei Zhou, Jingwan Tang, Beilei Guo, Hanjia Lyu, and Zhen Bai. 2022. Challenges and Design Opportunities in Data Analysis for ML-Empowered Scientific Inquiry - Insights from a Teacher Professional Development Study. In *Proceedings of the International Society of the Learning Sciences Annual Meeting*. 847–854.