# MindScratch: A Visual Programming Support Tool for Classroom Learning Based on Multimodal Generative AI

Yunnong Chen
College of Computer Science and
Technology, Zhejiang University
HangZhou, ZheJiang, China

Shuhong Xiao
College of Computer Science and
Technology, Zhejiang University
HangZhou, ZheJiang, China

Yaxuan Song
College of Computer Science and
Technology, Zhejiang University
HangZhou, ZheJiang, China

Zejian Li
School of Software Technology,
Zhejiang University
HangZhou, ZheJiang, China

Lingyun Sun
College of Computer Science and
Technology, Zhejiang University
HangZhou, ZheJiang, China

Liuqing Chen[*]
College of Computer Science and
Technology, Zhejiang University
HangZhou, ZheJiang, China
chenlq@zju.edu.cn

## Abstract

Programming has become an essential component of K-12 education and serves as a pathway for developing computational thinking skills. Given the complexity of programming and the advanced skills it requires, previous research has introduced user-friendly tools to support young learners. However, our interviews with six programming educators revealed that current tools often fail to reflect classroom learning objectives, offer flexible, high-quality guidance, and foster student creativity. This highlights the need for more adaptive and reflective tools. Therefore, we introduced MindScratch, a multimodal generative AI (GAI) powered visual programming support tool. MindScratch aims to balance structured classroom activities with free programming creation, supporting students in completing creative programming projects based on teacher-set learning objectives while also providing programming scaffolding. Our user study results indicate that, compared to the baseline, MindScratch more effectively helps students achieve high-quality projects aligned with learning objectives. It also enhances students' computational thinking skills and creative thinking. Overall, we believe that GAI-driven educational tools like MindScratch offer students a focused and engaging learning experience.

## Keywords

Computational Thinking, Generative AI, Mind Mapping, Programming Support Tool

## 1 Introduction

In the past two decades, computational thinking (CT) has become a focal point in educational research and practice [20, 56]. While CT and programming are distinct concepts, most educational practices regard learning programming as an effective way to cultivate CT skills [63]. This perspective underscores the importance of early exposure to CT in shaping children's futures, leading to a recent trend of teaching programming to young learners [32]. Accompanying this development is the emergence of programming languages designed specifically for young learners, with Scratch being a prime example [45]. Scratch is widely used for creative programming, with projects often involving games, stories, and animations [50]. However, educators face a trade-off between rigid knowledge transmission and engaging project-based learning when teaching programming. Effectively guiding learners to stay aligned with predefined learning objectives and providing personalized feedback during project-based learning presents a significant challenge.

Previous research has explored methods of using advanced programming learning tools to support children's creative programming in self-directed learning. These tools can be categorized into programming concept learning tools and block-based programming support tools. For instance, Visual StoryCoder [12] focuses on mastering abstract concepts in visual programming. CodeToon [60] is dedicated to transforming abstract code into engaging stories and comics. ChatScratch [5] supports creative programming by providing both creative stimulation and code guidance. However, without structured courses and instructional guidance, the potential for meaningful learning outcomes is significantly diminished [34, 41, 46].

To understand the specific needs of creative programming in the classroom and assess the effectiveness of current programming support tools in such contexts, we conducted a formative study with six programming educators. Educators evaluated the applicability of these tools in classroom settings from the perspective of structured course design, focusing on objectives, processes, and content [48]. The formative study revealed three key insights: (1) Strategies teachers use to ensure students achieve their learning goals during project-based learning. (2) Approaches for providing customized creative and programming support. (3) Methods to make creative resources more accessible to students while remaining manageable

for teachers. Additionally, we found that previous programming support tools lacked controllable guidance for students, hindering alignment with course learning objectives and delaying timely feedback from teachers. This issue may result in teachers spending excessive time addressing simple questions instead of offering constructive feedback promptly.

Based on these insights, this paper introduces MindScratch, a programming learning tool designed to support students in creating personalized, high-quality programming projects in classroom settings. Table 1 summarized the applicability of MindScratch and other programming learning tools in classroom settings across five key aspects: alignment with class objectives, task breakdown and progression, access to multimedia materials, controlled learning guidance, and support for teachers in providing timely feedback. To better help students achieve their learning objectives, teachers can set initial task themes and learning objectives in MindScratch. The tool receives these inputs and presents them in the form of a mind map. Leveraging the system prompt capabilities of large language models (LLMs), MindScratch retains the teacher-defined learning objectives throughout student interactions. Guided by an LLM-driven conversational agent, MindScratch provides students with controlled learning guidance rather than simply answering all of their questions.

To inspire students' creativity, MindScratch encourages the generation of multimodal creative materials—such as text, images, audio, and code—through conversations aligned with the classroom workflow, covering project planning, material creation, and code implementation. These materials are displayed as a mind map based on the classroom workflow structure. MindScratch helps students break down teacher-set learning goals into sub-tasks, manage the learning process, use generated materials for personalized visual effects, and complete programming projects with code suggestions. Notably, with MindScratch's support, teachers can shift their focus from managing numerous students' project progress to providing timely and constructive feedback. To develop an effective Scratch educational LLM, we collaborated with educators to create a specialized dataset and fine-tuned the GPT-3.5-turbo model to provide scaffolded guidance. Additionally, MindScratch employs a text-to-image generative model [39] and a text-to-audio generative model [16] to help students access high-quality creative materials.

To evaluate MindScratch's effectiveness in facilitating visual programming in the classroom, we conducted a within-subject study involving 24 fifth-grade students. The results indicated that, by using MindScratch, students more effectively achieved their learning objectives, showing significant improvements in CT skills and creativity. The extensive expansion of mind map nodes demonstrated deep engagement and valuable explorations by the students. Additionally, long-term monitoring of three students, combined with semi-structured interviews with six programming educators, provided insights into the future adoption and integration of AI assistants like MindScratch into AI-supported curricula.

In summary, the main contributions of this work are:

- Through a formative study, we identified the challenges of aligning objectives, managing processes, and delivering content in creative classrooms, along with the limitations of current programming support tools.

- We introduced MindScratch [1], a visual programming support tool that employs interactive mind maps powered by multimodal generative AI. This approach provides interactive creativity exploration, scaffolded code assistance, and the generation of creative materials, thereby enhancing students' learning outcomes in the creative classroom.
- We conducted a comparative study between MindScratch and Scratch, highlighting MindScratch's effectiveness in enhancing students' programming learning. We also provide insights and design considerations for future creative programming support tools.

## 2 Related Work

### 2.1 Generative AI in Computing Education

Generative artificial intelligence (GAI) has become a significant topic of debate, particularly due to its ability to produce high-quality artistic media across visual arts, concept art, music, literature, video, and animation. For instance, diffusion models can synthesize high-quality images [39], and large language models (LLMs) can generate coherent and impressive prose and poetry in various contexts [3].

In the context of computing education, GAI offers several advantages. One major benefit is the enrichment of learning resources. Sarsa et al. [52] demonstrated the effectiveness of OpenAI Codex in generating personalized code exercises and complex code explanations, thereby reducing teachers' preparatory workload. MacNeil et al. [33] further explored how GAI can provide code explanations by examining the relationship between different prompts and the resulting explanations. Another advantage lies in supporting the programming process, notably through the provision of exemplar solutions. These solutions offer students an easily understandable starting point [5, 26], enhancing their task performance. Additionally, students can learn from these examples by observing various problem-solving approaches, coding styles, and program designs [13, 22].

However, it is crucial to address the challenges and negative impacts of GAI in educational practice. Researchers have expressed considerable concern about the misuse of GAI, especially among novices introduced to programming in the era of such technology [10, 25, 69]. Chen et al. [6] highlighted potential harmful biases in generated code, including issues like racism and stereotypes. Security concerns are another significant issue. Moreover, instances of academic dishonesty have increased, with students turning to GAI to complete coding assignments [13, 14]. The variability of GAI-generated code makes such cheating more difficult to detect [57]. Over-reliance on GAI may also reduce students' critical thinking, as novices might adopt generated solutions without thoroughly engaging with the material, which is detrimental to learning processes based on reinforcement [2]. To mitigate these issues, we employ mind maps to visualize the project development process, collaboratively developing the mind map and programming project with students in a step-by-step approach.

---

[1]https://github.com/ArthurWish/MindScratch

**Table 1: Comparison of MindScratch with multiple programming learning tools based on educators' feedback, with a focus on their applicability in classroom programming instruction.**

| Targets | CodeToon [60] | Visual StoryCoder [12] | ChatScratch [5] | MindScratch |
|---|---|---|---|---|
| Reflect Class Objectives | ✓ | ✗ | ✗ | ✓ |
| Task Breakdown and Progression | ✗ | ✓ | ✗ | ✓ |
| Access to Multimedia Materials | ✓ | ✗ | ✓ | ✓ |
| Controlled Learning Guidance | ✗ | ✗ | ✗ | ✓ |
| Support Teachers to Provide Timely Feedback | ✗ | ✗ | ✗ | ✓ |

## 2.2 Learning Computational Thinking through Programming

To meet the growing demand for 21st-century skills, efforts are being made worldwide to integrate computational thinking (CT) into school education [20, 62]. While programming is not the only way to teach CT, its practical methods and strong connection to CT concepts have made it a key approach in classrooms [63]. Some researchers argue that programming should be seen as a teaching strategy rather than just a technical skill or set of coding techniques. Based on this view, integrating programming with classroom instruction can lead to more meaningful outcomes [35, 49].

Efforts are also being made to reduce teaching workloads by using AI systems to handle basic tasks. For example, VizProg [74] offers a system that helps teachers monitor students' coding progress in real time, allowing them to assist students who may be struggling. Tools like PuzzleMe [64] and VizPI [61] promote peer assessment in programming classes, encouraging collaboration and code evaluation among students. Additionally, many other tools guide students in areas such as variable declaration [17], coding style [36], and code correctness [58].

To support the widespread adoption of programming education in K-12, particularly in the early grades, significant efforts have been made. Due to challenges related to literacy demands and code conventions, block-based programming [28, 44, 45] has become a popular method for teaching CT. These approaches are often combined with specific tasks like robot programming [18, 47] or creative programming [5, 50] to meet learning goals and increase student engagement. To address the lack of collaboration in programming systems, Wang et al. [66] introduced a new tangible, collaborative programming system called Lighters. However, classroom teaching is a structured process where teachers carefully plan lessons for different periods, and previous studies have not adequately provided controlled learning guidance to support teaching. This study explores the use of mind maps as a visual tool to help teachers better manage students' progress and assist students in organizing project ideas and programming materials more effectively.

## 2.3 Assisting Novices in Learning Programming

Romero et al. [50] divided programming learning into five stages based on learner engagement. The first stage involves passive exposure to teacher-led explanations, videos, or tutorials, where learners are not yet involved in the creative process. This stage is supported by widely available open-source videos and tutorials. The second stage introduces step-by-step programming activities, allowing learners to follow instructions to achieve predefined goals. Platforms like Code.org [8] and Hour of Code [19] provide numerous projects for novices to gain hands-on experience. Recently, LLMs have also been used to generate tutorials for simple tasks [27, 43].

The third stage involves creative coding, where learners use programming tools to develop original works, integrating knowledge across multiple domains. Tools at this level typically offer creativity support and process control to balance task completion with divergent thinking [5, 11, 12, 65]. Building on this framework, our study investigates the potential of conversational agents powered by generative AI to provide structured guidance aligned with classroom workflows while delivering an engaging creative programming learning experience.

## 3 Formative Study

Our formative study aims to identify the challenges teachers and students face in creative programming classrooms, as the teaching is project-based, support-consuming, and content resource-consuming, which is special compared to other classroom settings.

### 3.1 Procedure

We recruited six educators (3 females, and 3 males) of age 26-39 (M=31, SD=4.44), each with over four years of experience in teaching programming to students and proficient in multiple programming languages, including Scratch, Python, and C. All educators are affiliated with a prominent programming education institution in China and have substantial experience in leading both small after-school courses (3-9 students) and larger school-based courses (over 20 students).

We conducted one-on-one discussions with each educator, introducing them to three representative programming support tools: CodeToon [60], Visual StoryCoder [12], and ChatScratch [5], all of which were presented at prominent conferences (e.g., CHI and UIST) in the past two years. CodeToon allows teachers to bridge the gap between abstract coding concepts and tangible learning outcomes by transforming symbolic text code into engaging stories and comics. This helps students visualize and better understand complex ideas. Visual StoryCoder offers a unique, block-based language that simplifies programming concepts, ideal for teachers aiming to introduce foundational programming skills. ChatScratch, utilizing generative artificial intelligence, enables students to actively engage in creating story-based projects in Scratch, fostering creativity and problem-solving skills in a supportive, AI-enhanced environment.

Before our discussions, none of the educators had any knowledge of these three applications. To further prevent the researchers' biases from influencing their attitudes, we facilitated their understanding of these applications in two ways. First, we presented video

materials of these works, which were uploaded to the ACM Digital Library to the educators. Second, the educators themselves experienced the usage of these tools. For tools that were not open-source, we conducted simulation using Wizard of Oz (WOz) techniques [15]. The simulations were based on the descriptions provided in the methods section of the respective papers and application usage procedures (if exist) documented in their supplementary materials.

## 3.2 Findings

Systematic programming teaching is typically anchored in a curriculum that is meticulously designed by teachers, encompassing the objectives, process, and content of each lesson. In this case, educators expect supporting tools to enhance student's learning experience and provide additional convenience while being adaptable to the existing teaching approach. As shown in Table 1, a summarized takeaway for our study is that these three representative tools faced challenges to be directly applied within classroom settings for creative programming. Based on feedback from educators, we report the main issues as below.

*3.2.1 Objectives.* Educators expect the support tools to align with the specific objectives they have designed, rather than allowing students to engage in unrestricted activities. This necessitates that teachers possess the ability to directly tailor and intervene with the support tools. For example, as P2 suggested, teachers should be able to *"specify the knowledge points that need to be addressed or the themes of programming tasks in the current lesson."* In this context, CodeToon stands out for its ability to effectively facilitate the exploration and understanding of designated code fragments. Conversely, Visual StoryCoder does not offer the functionality for educators to pre-define learning objectives, which means its pedagogical effectiveness is wholly dependent on the learner's self-directed engagement. Regarding ChatScratch, it accommodates the specification of programming themes, yet due to its focus on supporting creativity with AI-generated content, the control over programming project workflows is notably compromised.

*3.2.2 Process.* Classroom time is always limited, and teachers must wisely allocate learning time for each key objective, striking a balance between fostering students' divergent exploration and advancing the process. In this case, educators prioritize the system's ability to function as a personal supervisor for each student, managing the learning pace and keeping students focused. For instance, *"the system could provide a detailed breakdown of steps based on the task objectives and guide students through completing each item sequentially"*, as noted by P1. This ensures that students progress through their lessons efficiently, adhering to a structured timeframe, and keeps the class on track. In this regard, Visual StoryCoder offers process guidance through a conversational agent, clarifying the objectives of each step, thereby facilitating task progression. Conversely, CodeToon and ChatScratch allow for free exploration, which may lead students to focus on refining comic details or generating assets, diverting attention from their primary goal of programming.

Educators also highlighted key barriers that slow the class process, including aspects related to creativity, such as the conception of project roles and plots, as well as parts related to coding, such as the implementation of logic and data flow. Each student faces personalized issues, requiring teachers to address them specifically. *"Only after resolving everyone's confusion, the class can progress to the next step.",* as said by P1. In this case, they expected the system to provide support when critical confusion arose, pulling students out of any standstill. For CodeToon and Visual StoryCoder, the answer is negative. ChatScratch provides two-stage assistance in its code support process, including suggestions for code selection and code implementation. Nonetheless, educators believe that the generated code in this process remains difficult to comprehend, especially for beginners.

*3.2.3 Content.* In the field of programming education, project-based learning (PBL) has become a highly popular teaching method. When students are learning by solving real-world problems, they always need to utilize additional content resources. For instance, when constructing a simple website, students may require additional textual and decoration materials; while creating animations in Scratch could involve the introduction of character images and audio materials. Generally, these materials are prepared in advance by teachers, although students are sometimes dissatisfied with them due to a lack of choice and personalization. While allowing students to source their materials online can be too time-consuming and may lead to classroom disorder. Within this context, educators expect support tools to facilitate access to the additional materials needed for the projects. On this matter, CodeToon is capable of providing additional stories and comics, yet there is a convergence issue in terms of form; While Visual StoryCoder lacks this feature. Leveraging the artistic capabilities of generative AI, ChatScratch can assist students in obtaining the necessary materials, with its support primarily focused on images.

The quality of materials provided by the system is also a major concern for educators. They need to be accurate, non-toxic, and appropriate for students. In this regard, educators believe that CodeToon can meet their needs because it relies on established rules to parse code and utilizes confirmed materials to create stories and comics. Regarding Visual StoryCoder, educators express concerns due to its reliance on students' active storytelling and drawing. Without proper guidance, there is a risk that students might create or come across content that includes inappropriate language, imagery, or concepts. For ChatScratch, researchers suggest the need for additional mechanisms to review AI-generated content.

## 3.3 Design Goals

Building upon the insights from our investigations, we derived three principal design goals for our system to support students in creative programming learning in the classroom:

- **DG1: Support alignment between the project development process and learning objectives:** The system should ensure that students' project creation process stays aligned with the classroom learning objectives set by the teacher. Therefore, the system needs to display the teacher's requirements, such as the story elements and code blocks that must be included in the project. To ensure that students pay attention to and understand these key elements and code blocks, the system should highlight this information and provide necessary explanations.

- **DG2: Provide guidance and real-time support for programming learning at scale:** To reduce the pressure on teachers in maintaining students' progress alignment and allow them to focus more on providing constructive feedback, the system should offer project development support for students. For example, it could visualize the students' ideation process, guide them in breaking down tasks, and offer suggestions for code implementation.
- **DG3: Provide creative programming support based on classroom projects:** In creative programming classes, students can build upon project templates provided by the teacher. However, the limited programming resources offered by the teacher may not meet the diverse creative needs of students, and teachers may lack the time to help each student acquire custom images or audio materials. Therefore, the system should provide multimodal creative material support, leveraging the generative capabilities of GAI to help students realize personalized programming projects.

## 4  MindScratch

Based on the design objectives, we developed MindScratch, a GAI-driven visual programming learning system to help students achieve creative programming projects aligned with classroom goals. Our system uses a mind map to display the project objectives set by the teacher. Teachers can create an initial mind map in the system and input the learning objectives through an input box. For example, a teacher can set a theme like "Kitten Fishing" (red node, Fig. 1.(a.3)) and course objectives (such as learning conditions and loops, gray nodes, Fig. 1.(a.3)) within the mind map. The LLM in the system receives the learning objectives as a system prompt, which serves as the context for generating subsequent content. Through the LLM-powered conversational agent (Fig. 1.b, chat box), the system guides students in task breakdown, sparks creativity, and offers coding implementation suggestions based on the learning goals set by the teacher. The content created by students during their interaction with the agent is then presented in the mind map (Fig. 1.a, mind map). Finally, students complete creative programming tasks based on the mind map they created. For example, they can import character images from the mind map into the Scratch software, use generated audio to enhance the storytelling aspect of their programming, and implement program functionalities according to the programming logic and code suggestions on the mind map. The advantage of the mind map is that it reduces students' cognitive load by breaking down the process of simultaneously conceptualizing ideas and programming into two distinct steps: first brainstorming in the mind map, and then proceeding to program.

### 4.1  MindScratch's User Interface

Figure 1 illustrates the user interface of MindScratch, which is composed of three primary panels: an interactive mind map with a block palette, a chat box, and a drawing board (Figure 2). The mind map allows nodes to be added in two ways. Users can either right-click on a selected node to choose a generated node type (character, logic, or code), or they can follow the suggestions provided by the agent during conversations. This visual representation helps students organize their projects. It also aids teachers in monitoring

individual progress and aligning learning objectives. To create an engaging learning experience, we integrated a chat box where the agent provides guidance and suggestions to help students construct their projects through interactive conversations. Additionally, we developed a drawing board and a text-to-audio generation feature, offering students easy access to high-quality programming assets.

*4.1.1  Mind Map view.* Existing research indicates that mind maps are effective tools for learners to grasp complex concepts, fostering engagement in critical, higher-order, and reflective thinking [21, 59, 70]. We organize characters, images, and code concepts into nodes within the mind map to help students form a structured understanding. Inspired by construct-on-scaffold mind mapping [77], MindScratch uses three types of nodes to organize the programming project's mind map. These three types of nodes are distinguished by different colors: characters (green), logic (purple), and code (gray). We adopted a Scratch-style block palette (Figure 1.(a.4)) that allows students to select and add Scratch blocks.

Another advantage of using mind maps is that we can provide more explanatory information when students use MindScratch. An additional LLM is instructed to annotate the generated content and serve as edges between nodes in the graph (Figure 1.(a.5)). To avoid information overload, MindScratch only shows brief relationship at a time, and students can click the relationship to check more. We introduced two features to address educators' challenges in using MindScratch to meet their clearly defined classroom objectives: setting up the initial mind map and defining objectives. Our system allows educators to construct mind maps as a starting point and input learning objectives to control the system's generation. When students are creating mind maps, we instruct an underlying LLM to ensure the generated content revolves around the learning objectives set by teachers. When students add nodes, our system detects whether the added nodes are relevant to the learning objectives and highlights those nodes if they are relevant (Figure 1.(a.6)). After completing the mind map, students can engage in hands-on creative programming based on the theme tasks outlined in the mind map. For example, students can use the code blocks from the mind map to implement specific features and import materials from the mind map. If they encounter issues during the programming process, they can still seek help from MindScratch for code suggestions. It is important to note that at no point does MindScratch generate a complete code solution directly; instead, it offers logic suggestions and key code blocks.

*4.1.2  Dialogic Question-Answering Box.* Based on prior research [40, 51], we designed a dialogue interface to facilitate deeper student reflection and provide support during the project development process. The bottom of Figure 3 illustrates how an educator sets programming themes and prompts learning objectives in MindScratch at the beginning of the class. The LLMs in MindScratch use the established learning objectives as memory, incorporating more content related to these objectives in their generation, such as highlighting mind map nodes or providing additional explanations. The top of Figure 3 describes the interactive process between students and MindScratch. In the conversation process, the system captures students' needs and generates suggestions in a multiple-choice format. For example, the system guides students to think about the required characters by asking: *"What characters will your project*
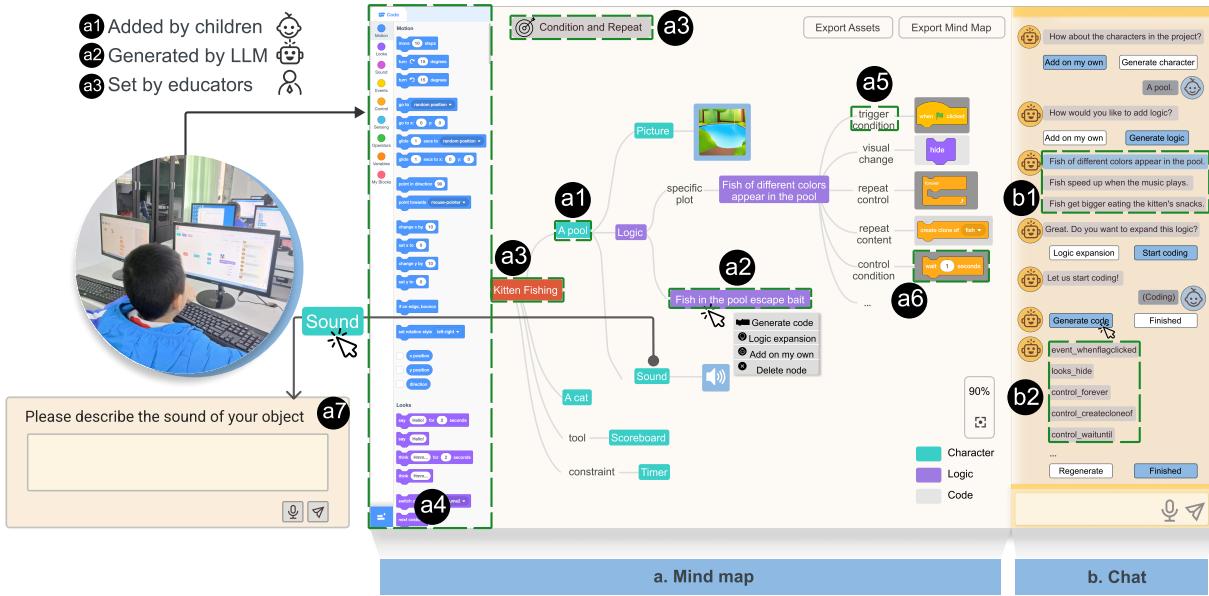
**Figure 1: MindScratch's User Interface. In the mind map (a) with block palette (a4), each node represents the collaborative creation of students, teachers, and AI (a1, a2, a3). The connections between nodes are annotated, and the highlighted nodes indicate their relevance to learning objectives (a5, a6). Within the dialogue box (b), the system provides structured guidance and real-time support through a Q&A mechanism (b1, b2). A floating text input box allows for the generation of audio materials through text editing (a7).**

*include?"* Students can then choose to manually add nodes or click the "character generation" button to select suggested characters for their mind maps. In the materials creation stage, if students click on a character node (marked in green), the system inquires, *"Would you like to add an image or sound to this interesting character?"* This interaction enables students to transition to an AI-driven drawing board or a text-to-audio input box, facilitating the creation of multimedia materials. During the code implementation stage, MindScratch provides scaffolding to encourage students to elaborate their programming ideas. Students can receive assistance by clicking on the "generate logic" or "generate code" button within the chat box. For instance, if the child provides a brief description of a character without much detail, the agent will ask a follow-up question, *"Very good! Would you like to add an action to the character?"* If the child is still confused, they can click the "generate logic" button in the chat window, and the system will generate various programming logic options for them to choose from (Figure 1.(b.1)). This feature particularly benefits beginners who may not fully grasp how algorithmic logic translates into actual Scratch code. For further clarification and assistance, the "generate code" button within the chat window allows students to generate Scratch code (Figure 1.(b.2)).

*4.1.3 AI-driven Drawing Board.* To facilitate the easy creation of elements, such as characters and scenes in the mind map, MindScratch offers an AI-driven drawing board for sketching and refining images (Figure 2). Students can utilize drawing tools to sketch, erase, select colors, and adjust line widths on the canvas. Given students' limited drawing skills, the system allows them to draw only basic sketches. By pressing the voice button (Figure 2.(c.1)) to input speech and clicking the "image polish" button (Figure 2.(c.2)), the

system automatically refines their sketches into high-quality image assets. Students can regenerate images iteratively until they are satisfied with the output. Students can save their sketches (Figure 2.(c.3)) or refined images, which will be displayed as thumbnails in the "created assets" column on the right. We also emphasize enabling students to create audio assets to enrich their projects. Once images are saved to the mind map, students can click the "sound generation" button and input text to generate interesting sound effects.

## 4.2 MindScratch's Backend Model

*4.2.1 Stage-based Dialogue Process.* MindScratch uses dialogue across three phases—project planning, material creation, and code implementation—to support students in developing creative programming projects. The dialogue is powered by an LLM. To ensure the LLM provides guidance based on learning objectives, we use the teacher's input as a system prompt to instruct the LLM to follow these learning goals when generating logic and code suggestions. As long and complex prompts tend to decrease the task performance of LLMs [3, 68], we used dedicated prompts for each stage instead of combining instructions for all phases in a single prompt. Figure 3 illustrates the stage-based dialog process for prompting LLMs. We instruct GPT-4 to support three primary tasks consistent with the classroom process: 1) Crafting questions that guide students in breaking down tasks, such as conceptualizing characters, actions, and events. 2) Formulating questions that help students create materials, including sketching and audio generation. 3) Offering help with programming logic or generating code, tailored to the students' input. To ensure the consistency of the content

generated by LLM, we use the mind map information as a memory to prompt the LLM at every stage.

*4.2.2 Relationship Annotation and Block Highlighting.* To enhance the interpretability of the content generated by LLMs, we employ few-shot prompting to instruct an additional GPT-3.5 in generating semantic relationships between nodes. We express it in triplets, such as <Theme, Relation, Character> and <Logic, Relation, Code>. However, annotating the relationships of all nodes might lead to the inclusion of less important relationships in the text, causing information overload for students. Therefore, we only label the edges connecting adjacent nodes generated by the LLM. To align the teacher's set learning objectives with the students' developed projects, we utilize the LLM to detect nodes related to the learning objectives within the mind map. The system then highlights these elements. For example, if the learning objective includes understanding conditional logic, the code blocks related to conditional logic will be highlighted. Drawing from prior work [23], we leverage saliency filters to manage the complexity of the mind map. We instruct GPT-4 to determine whether code blocks should be highlighted, designated as either *high* ($H) or *low* ($L). We only highlight code blocks that are highly relevant to the learning objectives, to reduce confusion among students.

*4.2.3 Multimodal Assets Generation.* MindScratch integrated advanced image generation technologies based on the Stable Diffusion model [39] and ControlNet [76] to refine students' doodles. ControlNet incorporates an extra step by taking an input image and employing the Canny edge detector to identify its outlines. The resulting image, highlighting these detected edges, is saved as a control map. This control map is then used as supplementary conditioning alongside the text prompt when feeding into the ControlNet model. It is then fed to Stable Diffusion as an extra conditioning together with the text prompt. The refined images are generated based on these two conditionings. MindScratch employed a text-to-audio model [16] based on the diffusion model to generate high-quality audio materials. It is worth noting that we instruct GPT-4 to transform students' inputs into prompts that align with both the image and audio generation models.

We conducted an image quality experiment to ensure that the images generated by MindScratch meet children's needs. The results show that the images created by MindScratch exhibit significantly better alignment with built-in materials compared to community-sourced assets. To quantify this, we used the Fréchet Inception Distance (FID) and Sliced Wasserstein Distance (SWD), which assess stylistic, textural, and structural similarities between image sets. Lower scores on both metrics indicate greater resemblance. As a baseline, the FID and SWD between built-in and community-sourced materials are 48.40 and 126.40, respectively. In contrast, the difference between MindScratch-generated assets and the built-in library is much smaller, with an FID of 8.12 and an SWD of 62.30. Additionally, we have taken measures to ensure that the generated content is safe and appropriate for students. To this end, we use a role-playing strategy with the LLM [55], instructing it to act as both an educator and a child's assistant. We also implement negative prompts [24] in the stable diffusion model, explicitly excluding keywords related to violence, horror, sex, crime, and discrimination, thereby safeguarding the content for a young audience.



**Figure 2: MindScratch's drawing board. It enables students to create, refine, and save image materials (c1, c2, c3).**

*4.2.4 Step-by-step Code Assistant with LLMs.* Figure 4 shows the logic-code assistance pipeline for MindScratch. We proposed a scaffolded code support process, which provides step-by-step logic expansion and code generation support. To maintain consistency with the support pipeline and reduce the time required for generation, we fine-tuned a large language model. As shown in Figure 4, we created 98 Scratch code samples from Scratch cards [2]. To fine-tune the LLM, we converted the Scratch code samples into pseudo-code form. We adopted chain-of-thought reasoning to control the LLM's output and enhance the quality of the response. Since the generated code is in pseudo-code form, we further visualize the generated code into the mind map. Specific details are described in the following subsections.

**LLM fine-tuning.** We employ the GPT-3.5-turbo model as the basis for our fine-tuning process. Initially, we curate a dataset specifically tailored to our objectives, consisting of examples that closely mimic the types of interactions and outputs we aim for the model to replicate. This dataset includes a mix of dialogues, coding logic descriptions, and code solutions. The data collection and processing process is shown at the top of Figure 4. In the Scratch community, Scratch cards emerge as essential introductory materials for programming. To obtain high-quality data for fine-tuning, we hired three educators, each with over four years of experience in teaching Scratch. Their task was to independently review all Scratch examples and select samples of educational significance. After an hour of review, educators selected 98 representative code examples. Subsequently, two researchers compiled these code examples into the fine-tuning dataset in the chain-of-thought style. BLEU and F1-score metrics are added to evaluate the performance of the LLM. Using 30 test samples based on Scratch cards, the fine-tuned LLM achieved a BLEU score of 32.1 and an F1 score of 80.9, compared to the vanilla LLM's BLEU score of 29.5 and F1 score of 73.2.

**Prompt Engineering.** Table 2 displays examples of prompt engineering for logic and code generation. A concrete example aids LLMs in eliciting the specific knowledge and abstractions required to complete a task [67, 72]. To ensure that the fine-tuned LLM adheres to our code assistance workflow, we transformed the collected code examples into logic-code pairs and formulated

---

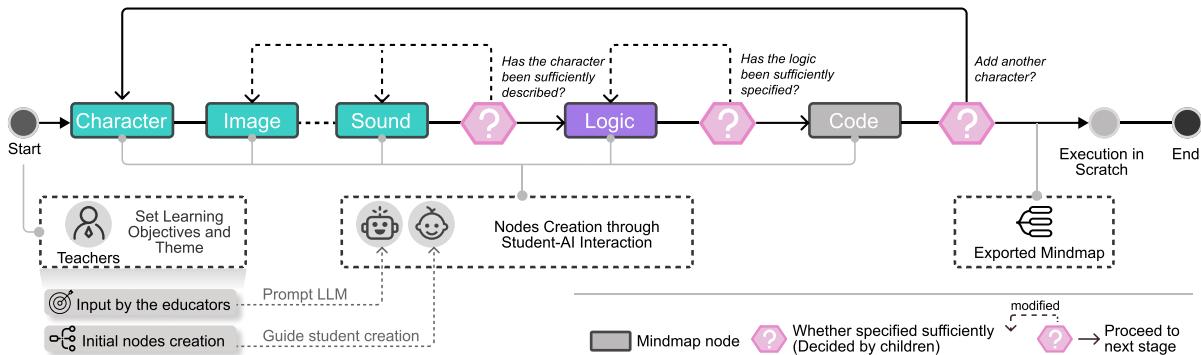[2]https://resources.scratch.mit.edu/www/cards/en/scratch-cards-all.pdf

**Figure 3: MindScratch user interaction process: Teachers set learning objectives and programming themes to prompt the LLM and guide student creation. Students collaborate with MindScratch to create characters, generate images and audio, and add them to the mind map. In the coding phase, students use LLM-provided logic and code blocks for implementation. Finally, students conduct hands-on Scratch programming and can ask MindScratch for help if needed.**
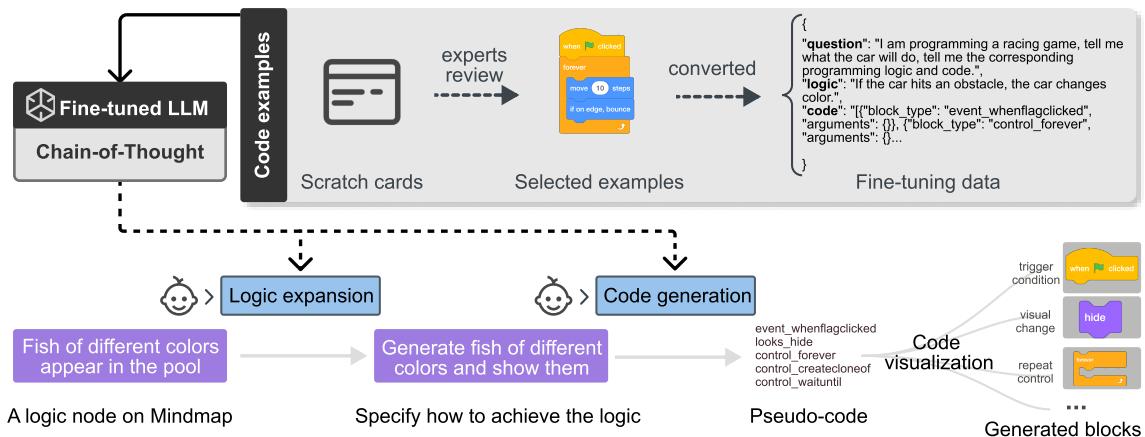


**Figure 4: Overview of the logic-code assistance pipeline for MindScratch.**

the instruction in a chain-of-thought style. During the testing by researchers, it was discovered that although LLMs can generate correct pseudo-code, there remain issues with parsing, especially with nested relationships, which leads to decreased match performance and, consequently, an inability to visualize. We instruct the fine-tuned LLM to generate an abstract syntax tree (AST), which preserves the hierarchical structure of the code, thereby facilitating more precise alignment.

**Visualization of Generated Code.** To convert the pseudo-code generated by the LLM into visual Scratch code for easier understanding by students, we collected screenshots and class IDs of all Scratch code blocks from the Scratch website. We use the edit distance to determine how far the pseudo-code block generated by the LLM is from the correct code block's class ID. We parse all the pseudo-code blocks from the AST generated by the LLM, then calculate the edit distance with all class IDs, and assign the class ID with the smallest distance as the generated block ID. For example, if the LLM outputs the pseudo-code "sensing touching_object Cat", it can be matched to the corresponding Scratch block ID "sensing_touchingobject". This matching approach also helps to minimize errors caused by hallucinations produced by LLMs [71]. We then use the matched

class ID to index the image of the code block, thereby visualizing the code block the students want, helping them correspond more easily with Scratch.

### 4.3 Implementation

MindScratch's mind maps are implemented using JSMind [3], and the user interface is built with Vue.js [4]. The backend framework was developed using Flask [5] in Python. We have utilized OpenAI's Whisper model to generate audio outputs for the conversational agent, and OpenAI's speech recognition API for transcribing speech [6]. To minimize the interference of simultaneous voices in the classroom, we also designed a typing interface. To run the underlying LLMs, we used OpenAI's ChatCompletion API [7]. After extensive testing, we chose the state-of-the-art GPT-4-0613 model to generate responses and translate students' prompts to the image and audio generation model. In terms of safety, we incorporated a moderation

---

[3] https://github.com/hizzgdev/jsmind
[4] https://github.com/vuejs/vue
[5] https://github.com/pallets/flask
[6] https://openai.com/research/whisper
[7] https://platform.openai.com/docs/guides/gpt/chat-completions-api

layer [8] to ensure outputs were devoid of sexual content, hate speech, harassment, violence, or self-harm. The final version of prompts used in our system is displayed in our supplementary materials.

## 5 EXPERIMENT

To investigate MindScratch's impact on students' creative programming process and outcome, we conducted a within-subject study, using the original Scratch platform as a baseline. Twenty-four fifth-grade elementary school students participated in our experiment. Each student used both systems, executing theme-based creative visual programming tasks in two separate sessions. Primarily, we aimed to address two research questions through our experiments:

- RQ1: How does MindScratch facilitate students' learning objectives in creative programming courses, reflecting the knowledge points and key tasks established by teachers?
- RQ2: How does MindScratch enhance the students' project code quality and computational thinking skills in creative programming courses?
- RQ3: How does MindScratch provide creative support for students in creative programming courses?
- RQ4: What are the perspectives of educators regarding learner-focused AI assistants like MindScratch in terms of its integration into the curriculum, recommendations for improvement, and effective pedagogy?

### 5.1 Participants

This study took place in a primary school classroom in Shenzhen, China, during the winter of 2023, involving 24 fifth-grade students. The participants were 24 fifth-grade students. All participants had already completed a semester-long Scratch course. All participants had a basic understanding of how to use Scratch's fundamental code blocks, as well as basic software operations such as uploading and editing images, and uploading and editing audio. However, they were not familiar with advanced CT skills, such as nested condition loops and the broadcast mechanism. To prevent the potential influence of different teachers on the outcome of the study, all students were taught by the same teacher using the same technical equipment regardless of which group they belonged to. The classes were taught by a teacher with 6 years of experience in teaching Scratch courses. The participants' ages ranged from 10 to 11 years. Each student engaged with two systems, Scratch and MindScratch, completing two theme-based creative programming tasks across separate sessions. All participants were native Mandarin speakers. As compensation for participating in the experiment, each student received a gift worth $15.

### 5.2 Procedure

As a within-subject study, each participant was required to undertake two theme-based creative programming sessions. To avoid potential time-related effects and carryover effects, we introduced a one-week interval between sessions. Moreover, we employed a counterbalanced design [9], alternating between combinations of tools (MindScratch and Scratch) and themes (A or B) to ensure the fairness of the outcomes. Inspired by previous studies and the

insights of educators [38, 75], we and the teacher designed two creative programming tasks: A) "Cat and Mouse" and B) "Kitten Fishing". Each session lasted about 80 minutes, with 20 minutes dedicated specifically to familiarizing with the MindScratch system and free exploration. Participants then had 10 minutes to ask questions, which were answered by researchers. The teacher spent 20 minutes explaining the key points of the programming task. The last 30 minutes were focused on the theme-based creative programming task. The specific details of students' theme-based creative programming tasks can be found in Appendix A. MindScratch and Scratch were both configured with a Chinese interface according to the participants' preferences. During the experiment, researchers were only allowed to assist the students in the event of unexpected technical difficulties. After each session, we collected screen recordings through capture software. Additionally, we collected the project files created by the students in sb3 format. Subsequently, the participants were asked to complete a creativity support index questionnaire [7] to evaluate how much effective creative support their process received. After completing these two parts, we conducted brief semi-structured interviews with the participants to gain deeper insights into their user experience with MindScratch for class-based creative programming tasks.

We also explored the long-term effects of MindScratch on the development of computational thinking in students. Three students from the class participated in an experiment during the winter vacation. We invited the teacher to use our tools to conduct a four-week one-on-three course, with one lesson per week. Each lesson involved programming one project using MindScratch. Before and after the course started, we conducted pre- and post-tests on them using a CT skills survey [29]. After the course ended, we asked them for their feedback and opinions on how the tool supported their programming in the classroom.

To further understand how educators utilize MindScratch in creative programming courses, we conducted semi-structured interviews with six educators in graphical programming. The interviews began with exploring the educators' backgrounds and their current challenges and strategies, particularly regarding students' creative exploration and programming implementation using LLM-based tools. Subsequently, we introduced the creative programming support assistant, discussed its capabilities, and the insights gathered from its long-term deployment, as summarized in Sections 6.1-6.3. The discussion then shifted to educators' perceptions of our tool: what they liked and disliked, their pedagogical and ethical considerations for using it, their interest and requirements for integrating it into their curriculum, and how they view its relationship with tools like ChatGPT. Each interview was conducted on Tencent Meeting and lasted about 1 hour.

### 5.3 Measurements

In this section, we describe the evaluation process, including the six categories of data collected for assessment. Table 3 displays the metrics used for quantitative assessment. Interviews based on artifacts were used for qualitative insights. Any data requiring subjective evaluation were independently assessed by multiple evaluators, with their agreement documented.

---

[8]https://platform.openai.com/docs/api-reference/moderations

**Table 2: Examples of prompt engineering for code generation.**

| PROMPT TYPE | INSTANTIATION |
|---|---|
| **Example of logic generation** | |
| **LLM instruction** | *Your task is to generate a programming logic node based on the provided mind map and corresponding question. The generated programming logic must not repeat the content of the mind map; the description must be concise and suitable for a child's level. The response should be in JSON format. Below are some examples.* |
| **The input example in LLM** | *Mind map: { "racing game": { "car": { "function": ["Keep the car moving", "Adjust the car's direction"] }, "track": { "function": ["Touching the track causes the track to change color"] } } }* |
| **User's input** | *Car.* |
| **The output example in LLM** | *{"logic": Adjust the car's direction.}* |
| **Example of code generation** | |
| **LLM instruction** | *Please provide well-structured Scratch code based on the programming logic entered by the user. Ensure that the generated code blocks can be correctly parsed by Scratch software. The response should be in JSON format. Below are some examples.* |
| **The input example in LLM** | *Here is an example: { "code": [ {"block_type": "event_whenflagclicked", "arguments": {}}, {"block_type": "control_forever", "arguments": {}}, {"block_type": "control_if", "arguments": { "condition": { "block_type": "sensing_touchingobject", "arguments": { "object": "obstacle" }}, "block_type": "looks_seteffectto", "arguments": { "effect": "color", "value": 100 }}, "block_type": "operator_subtract", "arguments": { "NUM1": { "block_type": "sensing_of", "arguments": { "property": "score", "object": "Stage" }}, "NUM2": 1 }}]} ] }* |
| **User's input** | *Keep the car moving.* |
| **The output example in LLM** | *[{"block_type": "event_whenflagclicked", "arguments": {}}, {"block_type": "control_forever", "arguments": {}}, {"block_type": "motion_movesteps", "arguments": {"steps": 10}}]* |

**Table 3: Summary of the collected data and the measures used in the experiments.**

| Data | Evaluation Metrics | Description |
|---|---|---|
| Code Quality Scores | Dr.Scratch Rubric | Assess code quality by quantifying seven computational thinking dimensions [37]. |
| Expert Ratings | Expert Ratings on Projects | Evaluation by Scratch educators to determine the logic consistency, quality, creativity, and originality of students' projects [1]. |
| Mind Map Node Count | Mind Map Richness | Count of graph nodes in projects to quantify project richness [30]. |
| Creativity Support Index Questionnaires | Creativity Support Index | Measure the usability and effectiveness of a system in enhancing creative programming tasks [7]. |
| CT Skills Pre and Post-Tests | CT Skills Survey | Measure the long-term potential impact of the system on the cultivation of computational thinking [29]. |
| Artifact-based Interview | Semi-Structured Interview | Gather insights on the creation of students' mind maps, classroom progress, and feedback. [42]. |

*5.3.1 Dr.Scratch Rubric Scores.* We adopted the Dr.Scratch scoring criteria [37] to assess the code quality of students' class-based programming projects. As a widely applied metric for evaluating Scratch code, it assesses code quality by quantifying seven computational thinking dimensions: abstraction, parallelism, logic, synchronization, flow control, interactivity, and data representation. The score for each dimension ranges from 0 to 3, representing three levels: basic, developing, and master. Essentially, higher code quality scores indicate that students have more extensively utilized and mastered CT skills during the project development process [12].

*5.3.2 Expert Ratings.* For this evaluation, three Scratch education experts were invited to assess the students' projects independently. The experts used a 5-point Likert scale based on four criteria adapting from [1] for evaluating the projects: (1) Consistency: indicating the extent to which a student's projects reflect the classroom learning objectives; (2) Quality: reflecting the intrinsic properties of the project, including completeness, clarity, detail, and harmony; (3) Originality: indicating the extent to which the project reflects the student's creation rather than being derived from existing materials; (4) Creativity: indicating innovative content expression and

rich imagination in the project's materials. The evaluation results showed a high degree of consensus among the experts, with an overall intraclass correlation coefficient (ICC) of 0.78 ($p < 0.05$).

*5.3.3 Creativity Support Index.* The creativity support index (CSI) [7] measures users' perceptions of tool usability and support for creative tasks. This survey consists of 12 questions distributed across six themes. We employ a 5-point Likert scale to collect feedback and ensure each student is in an undisturbed environment, allowing us to gather firsthand feedback from the students and gain deeper insights into their experiences.

*5.3.4 Mind map Node Count.* In the classroom, teachers' ideas often influence students' project planning. Inspired by [30], we adopted graph node count to reflect the richness of the mind maps used for programming. In our evaluation, the number of graphical nodes is counted as the sum of character nodes (project planning and multimodal materials) and programming nodes (logic and code blocks).

*5.3.5 CT Skills Survey.* To investigate the long-term impact of MindScratch on students' CT skills, we employ a 5-point Likert

scale revised CT skills survey proposed by [29], which includes creativity (4 items), cooperativity (4 items), critical thinking (4 items), algorithmic thinking (4 items), and problem-solving (4 items). Before the four-week course began, we conducted a pre-test on the students using the CT skills survey, and a post-test was administered after the course concluded.

*5.3.6 Artifact-based Interview.* We conducted interviews using two artifacts: the mind maps and the projects created. This allows us to gain a deeper understanding of whether MindScratch effectively supports learning objectives and creative programming. We structured our interview around three themes to understand the students' projects created around the classroom objectives provided by the teacher (questions 1-2), their classroom process (questions 3-4), and their views on MindScratch (questions 5-7):

(1) Tell me about your projects.
(2) Did you follow the classroom objectives? Was your idea realized?
(3) We noticed that during [specific time/event], you exhibited [specific behavior]. What problems did you encounter at that time?
(4) How did you address these challenges?
(5) Would you be willing to use MindScratch in class?
(6) Which features of MindScratch do you find useful?
(7) Do you think using MindScratch would make programming easier for you?

## 5.4 Data Analysis

We implemented a counterbalance design in our experiment to mitigate the effect of order and theme on the programming outcomes. The paired t-test was used to investigate the differences between Scratch and MindScratch. A Bonferroni correction is applied to mitigate the risk of Type I errors associated with multiple comparisons. All the data and analysis codes are included, for details, please see our supplementary materials.

## 6 results

## 6.1 RQ1: Impact on Achieving Learning Objectives

Figure 5 shows the distribution of scores by experts. In the same class time, students using MindScratch outperformed those using Scratch in terms of consistency, originality, and creativity. All 24 students who used MindScratch completed the teacher's predefined tasks and learning objectives. For instance, the objectives set by the teacher included making a cat follow the mouse cursor, having multiple mice appear, and each mouse moving continuously, among others. However, only 13 students using Scratch achieved the learning objectives. Although many projects created by students using Scratch contained interesting interactions, it is clear that these projects did not reflect the classroom learning objectives. Without timely intervention from the teacher, these projects would lack completeness. In particular, Table 4 shows that MindScratch plays an important role in promoting the consistency between programming projects and learning objectives ($t(23) = 5.32, p < 0.01^{**}$). We found that experts appreciated the projects created using MindScratch. Significant differences

were noted in originality ($t(23) = 6.86, p < 0.01^{**}$) and creativity ($t(23) = 7.60, p < 0.01^{**}$) between MindScratch and Scratch. For instance, P15 used MindScratch to craft a narrative about a fortunate mouse and a less fortunate cat. In the story, when the owner is away, the mouse ingeniously disguises itself as a lion to frighten the cat away. This story not only reflects the classroom learning objectives but also includes rich and exquisite assets, showcasing the student's creativity effectively. Besides, in terms of quality, we observed matched means (4.10 for MindScratch and 4.04 for Scratch) and standard deviations (0.59 for MindScratch and 0.69 for Scratch), which suggested that projects created with MindScratch in the classroom have the same level of completeness as those made with Scratch.
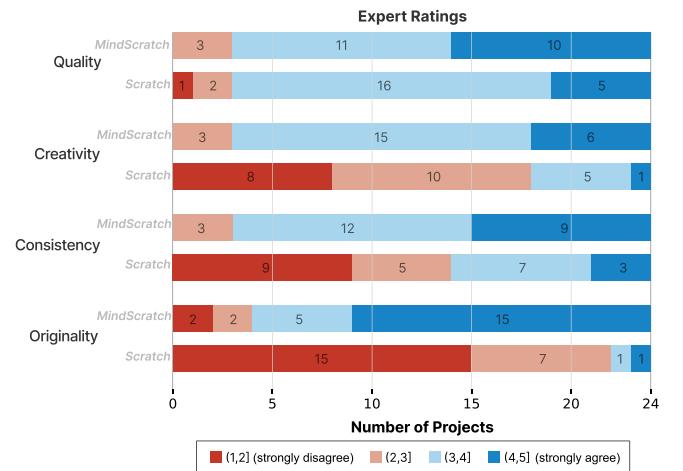


**Figure 5: The score distribution between MindScratch and Scratch. Note that higher values indicate positive feedback, and vice versa.**

## 6.2 RQ2: Effectiveness on Code Quality and CT Skills

Table 5 shows the results of the Dr.Scratch rubric, which evaluates the quality of code by assessing the level of CT in students' projects. For the total score, the mean for MindScratch was 14.17 ($SD = 4.14$)), while that for Scratch was 9.96 ($SD = 4.46$). Within the framework of Dr.Scratch, this progression marks the elevation of students' projects from a "basic" to a "master" level. Based on the paired-samples t-test, we observed an improvement in six dimensions: abstraction ($t(23) = 3.11, p < 0.01^{**}$), parallelism ($t(23) = 3.42, p = 0.011^*$), logical ($t(23) = 5.79, p < 0.01^{**}$), synchronization ($t(23) = 3.50, p = 0.016^*$), interactivity ($t(23) = 3.98, p < 0.01^{**}$), and data ($t(23) = 3.02, p = 0.049^*$), except flow control ($t(23) = 2.07, p = 0.399$). These results suggest that compared to the baseline tool, MindScratch led to a better programming outcome in terms of code quality and logic organization. In the interview, ten students indicated that the interactive mind map and the scaffolded logic-code generation contributed to their programming output. *"During the use of MindScratch, the code suggestions it gave me were helpful when I didn't know which code blocks to use to implement the logic"* (P14). *"With MindScratch, I could figure out*

**Table 4: Comparative Evaluation of MindScratch vs. Scratch Across Expert Ratings.**

| Metric | MindScratch | | Scratch | | Paired-t test | |
|---|---|---|---|---|---|---|
| | Mean | SD | Mean | SD | t | p |
| Originality | 4.38 | 0.97 | 2.48 | 0.88 | 6.86 | 0.000** |
| Consistency | 4.25 | 0.67 | 3.0 | 1.32 | 5.32 | 0.000** |
| Creativity | 4.13 | 0.61 | 2.91 | 0.93 | 7.60 | 0.000** |
| Quality | 4.10 | 0.59 | 4.04 | 0.69 | 1.30 | 0.83 |

**Notes:**
1. A Bonferroni correction is applied to mitigate the risk of Type I errors associated with multiple comparisons.
2. ** denotes $p < 0.01$ and * denotes $p < 0.05$.

**Table 5: Comparative Evaluation of MindScratch vs. Scratch Across Dr. Scratch Rubrics.**

| Metric | MindScratch | | Scratch | | Paired-t test | |
|---|---|---|---|---|---|---|
| | Mean | SD | Mean | SD | t | p |
| Abstraction | 2.17 | 0.56 | 1.71 | 0.55 | 3.11 | 0.039* |
| Parallelism | 2.19 | 0.51 | 1.58 | 0.50 | 3.42 | 0.011* |
| Logical | 2.04 | 0.81 | 1.13 | 0.34 | 5.79 | 0.000** |
| Synchronization | 2.08 | 0.95 | 1.46 | 0.66 | 3.50 | 0.016* |
| Flow Control | 1.92 | 0.50 | 1.63 | 0.49 | 2.07 | 0.399 |
| Interactivity | 2.13 | 0.61 | 1.54 | 0.51 | 3.98 | 0.005** |
| Data | 2.04 | 0.20 | 0.91 | 1.41 | 3.02 | 0.049* |
| Total Score | 14.17 | 4.14 | 9.96 | 4.46 | 6.44 | 0.000** |

**Notes:**
1. A Bonferroni correction is applied to mitigate the risk of Type I errors associated with multiple comparisons.
2. ** denotes $p < 0.01$ and * denotes $p < 0.05$.

the logic for implementing code without needing help from a teacher" (P2). We do not observe a significant effect on the difference in flow control. According to students' feedback, we suppose that this can be attributed to the understanding and application of complex conditional logic, which requires a higher level of computational thinking [53].

Figure 6 illustrates the pre-test and post-test results using the CT Skills Survey. Each student's CT skills improved, with P3's pre-test mean at 3.4 and post-test mean at 3.9, P4's pre-test mean at 3.9 and post-test mean at 4.05, and P21's pre-test mean at 3.65 and post-test mean at 3.95. Overall, there was an improvement in students' creativity, cooperativity, algorithmic thinking, and problem-solving. All three students mentioned they could more easily understand the relationships between code blocks (algorithmic thinking) and could design a complex programming plan (problem-solving). P3 stated, *"After using MindScratch, I learned **how to use the code from previous projects** to solve current problems."* P21 mentioned, *"I now like to create a mind map to help me **conceptualize programming**; with a mind map, I can complete a project faster."* We can observe that the creativity of the three students and the collaboration skills of the two students have been improved. P4 commented, *"When creating programming projects, I want to use a variety of materials (images, audio). But each time, the materials provided by the teacher are limited. MindScratch allows me to access the materials I want more quickly."* MindScratch does not show significance in the development of critical thinking. In future work, we plan to incorporate strategies from [73] for critical thinking training to nurture this essential skill.

## 6.3 RQ3: Effectiveness on Creativity Support

Another objective of our study focuses on how MindScratch supports students in enhancing their creativity. First, We employed the metric of mind map node count to evaluate the students' efforts in classroom-based programming. Based on the learning objectives, the teacher created two complete mind maps as the baseline based on the two themes. The average count in MindScratch was 52.45 (SD = 6.93), while the baseline average was 35. This fully demonstrates that students can express their creativity with MindScratch based on the learning objective.

Second, We also report on students' subjective perceptions of the creative support they received during the programming process. From Figure 7, we can see that students have a positive view of both systems and MindScratch received higher scores across all six dimensions. A paired-sample t-test shown in Table 6 suggests that MindScratch showed significant improvement in areas such as collaboration ($t(23) = 3.88, p < 0.01^{**}$), exploration ($t(23) = 3.04, p = 0.035^{**}$), expressiveness ($t(23) = 3.52, p = 0.011^*$), immersion ($t(23) = 2.92, p < 0.01^{**}$) and results worth effort ($t(23) = 4.11, p < 0.01^{**}$) compared to Scratch. This indicates that MindScratch provided better creative support to students while engaging in classroom project programming, enhancing their learning experience. Statistically, 12 out of 24 students mentioned the contribution of character generation, 21 highlighted the impact of image polish, and 8 acknowledged the audio generation. On average, the students used character generation 5.32 times (SD = 1.56), image polish 6.27 times (SD = 2.74), and audio generation 2.34 times (SD = 0.38). P23 stated, *"It made classroom learning more interesting for me. I could create funny images and sounds on my own, which were different from what the teacher provided. When I was unsure about how to draw the characters, it also offered help."* In conclusion, by interactively creating mind maps (collaboration), students' efforts in creative tasks become more meaningful (results worth the effort). The visualization based on mind maps aids in exploring the details of projects (exploration), while multimodal asset creation enables students to effectively express their inspiration (expressiveness). Scaffolded logic-code support offers an efficient pipeline that
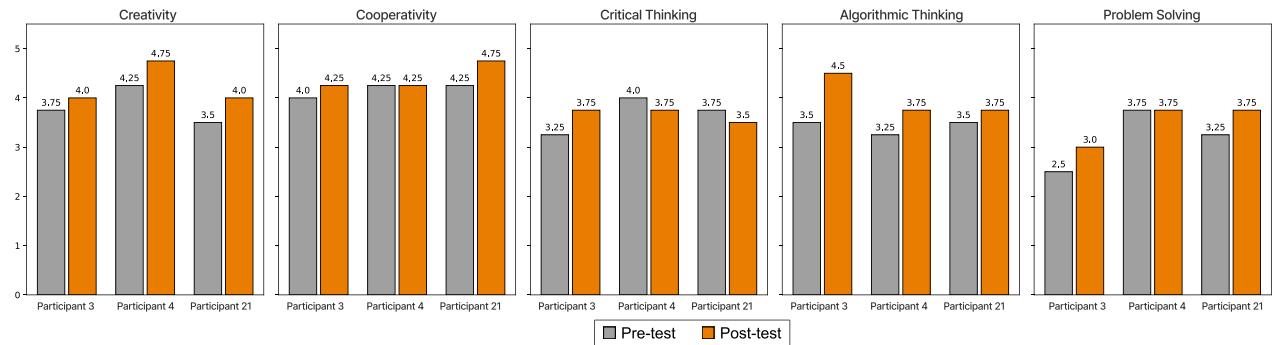
**Figure 6: The bar plots illustrate the pre-test and post-test results of three students across five sub-dimensions of the CT skills survey.**

**Table 6: Comparative Evaluation of MindScratch vs. Scratch Across Creativity Support Index.**

| Metric | MindScratch | | Scratch | | Paired-t test | |
|---|---|---|---|---|---|---|
| | Mean | SD | Mean | SD | t | p |
| Collaboration | 3.98 | 0.99 | 2.83 | 0.89 | 3.88 | 0.005** |
| Enjoyment | 3.73 | 0.69 | 3.35 | 0.80 | 1.72 | 0.59 |
| Exploration | 3.90 | 0.97 | 2.94 | 0.97 | 3.04 | 0.035* |
| Expressiveness | 3.94 | 0.86 | 3.02 | 1.02 | 3.52 | 0.011* |
| Immersion | 3.67 | 0.93 | 2.83 | 0.88 | 2.92 | 0.005** |
| Results Worth Effort | 4.19 | 0.83 | 3.23 | 0.78 | 4.11 | 0.003** |

**Notes:**
1. A Bonferroni correction is applied to mitigate the risk of Type I errors associated with multiple comparisons.
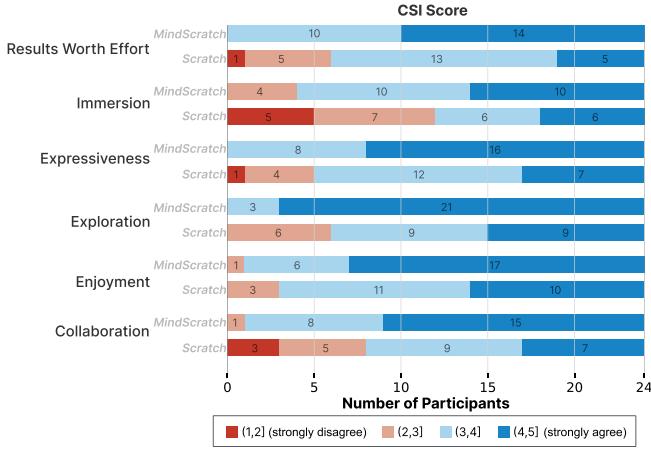2. ** denotes $p < 0.01$ and * denotes $p < 0.05$.



**Figure 7: The score distribution between MindScratch and Scratch. Note that higher values indicate positive feedback, and vice versa.**

allows for conceptualization and coding to be completed in one go. This enables students to focus more on creative programming (immersion).

## 6.4 RQ4: Educator's Perspectives

To gain further insights into how educators would use MindScratch in their visual programming classes, we conducted semi-structured interviews with six teachers (T1 - T6). Their professional opinions

provided more evidence on the appropriateness and effectiveness of using MindScratch in the classroom.

*6.4.1 General Impressions.* Educators generally have a positive impression of MindScratch. T1 emphasized the pedagogical approach of MindScratch, describing it as offering a "mind map creation method with AI collaboration," which is particularly helpful for students with unclear thoughts. Similarly, T3 mentioned that MindScratch is a more controllable way for elementary students to use large models," different from tools like ChatGPT that directly provide answers, instead, it guides students step-by-step to find answers. Furthermore, T5 envisioned that MindScratch could greatly assist teachers in dynamic creative classroom settings, helping teachers more easily handle spontaneous student issues, provide tailored feedback, and support teachers in quickly identifying students' knowledge gaps.

*6.4.2 Perceptions on Mind-map Usage.* Most educators appreciate the design of the interactive mind map, especially how it provides explanations after students select nodes and highlight code blocks related to the set educational objectives. T6 liked that it facilitates creative exploration but helps students avoid excessive divergence. T2 and T3 mentioned how it reduces cognitive load by visualizing the overall logic of projects. T4 stated that "providing logic-to-code guidance" aids in students' metacognitive skills. However, T1 expressed concerns about the lack of support for advanced algorithms, as the generated code blocks struggle to reflect the characteristics of data structures.

*6.4.3 Perceptions on External Materials Generation.* Most educators believe that using GAI to generate programming materials can

enhance students' learning interests and stimulate their creativity. Moreover, T2 stated that having students create their materials reduces the teaching burden on teachers to support differentiated student needs. T4 believes that before using MindScratch, creative projects merely reflected the teacher's creativity, rather than projects meaningful to the students personally. T3 and T5 discussed how to enhance the relationship between material creation and classroom themes by combining teacher-set mind maps. However, T1 expressed slight concerns, particularly that students focus on drawing or debugging prompts to generate content, neglecting subsequent programming training.

*6.4.4 Concerns about Incorrect Responses and Misuses.* Despite positive impressions, educators expressed various concerns. After viewing recordings of our system in use, T1 pointed out the risks of incorrect answers, especially for students with weaker foundations. T3 was concerned about students "believing anything the large language model says." As a solution, T1 and T3 suggested including systematic courses and quizzes before students use MindScratch. T5 asked to improve the model's accuracy by having students identify incorrect answers and suggested that the system should highlight potential uncertainties, prompting students to reflect, such as "Are you sure this code block will achieve the correct logic?" Additionally, T1, T2, and T3 expressed concerns about students misusing MindScratch and suggested limiting the number of daily generations as a potential solution.

*6.4.5 Student Monitoring Dashboard.* Another popular topic was the need for a teacher dashboard to monitor student interactions with the system and track their progress. T4 emphasized that by monitoring the questions students ask and the types of answers generated, educators can assess students' grasp of computational thinking, which can be obtained by the system generating reports on common issues. These insights could prompt interventions in the next class to provide better examples. However, accessing students' private data may be restricted. This raises a question about the balance between personalization and privacy. Instead of directly reporting student problems, the system should deal with potential privacy leaks first. T1 mentioned that students should not feel monitored while using the system.

## 7 Discussion

In this section, we reflect on our findings from the development and evaluation of MindScratch. We also discuss future opportunities for research exploring visual programming support tools.

### 7.1 Design Considerations

*7.1.1 Enhancing generated Answer Reliability and Providing Reliable Sources.* Based on our interviews, educators are still concerned that students may believe everything the LLM says, even if the answers may be incorrect. Although MindScratch prompts the LLM to annotate the relationships between nodes to improve the interpretability of the generated content, potential errors are still unavoidable. Therefore, MindScratch should incorporate advanced Retrieval-Augmented Generation [31] to enhance the reliability of its responses, while also listing the sources referenced when providing answers. A high-quality programming knowledge base can greatly improve the quality and reliability of the LLM's generated responses.

*7.1.2 Balancing Support and Self-individual Effort.* Drawing from the theory of scaffolding [54], MindScratch adopts a question-answering format to encourage users to exert effort before receiving support and feedback. For example, after children create digital materials, the system poses questions to guide them in thinking about programming logic. During the programming phase, MindScratch provides a Scratch-style block palette to encourage children to think about code implementation on their own. Our user study indicates that children actively expand on the code suggestions provided by the system during the programming, reflecting their interests. However, the system does not take into account the cognitive and capability development of children, and we suggest that future visual programming assistance systems should provide support and feedback according to the user's programming level. Moreover, MindScratch could incorporate gamification features (e.g., leaderboards, and online project sharing) to encourage users to enrich their projects for better performance.

*7.1.3 Providing Teachers with Student Monitoring Dashboard.* MindScratch visualizes students' thought processes during project creation through mind mapping, helping teachers monitor student progress and provide targeted feedback. However, educators believe that collecting additional process data (such as question data) would give a more comprehensive understanding of students' knowledge levels. MindScratch should offer a teacher dashboard to help educators track student progress and common questions. In future versions, MindScratch could further enhance collaboration between teachers and AI. For example, the system could allow teachers to set LLM prompts to provide assistance with varying levels of granularity based on classroom progress. Additionally, teachers could adjust the LLM's knowledge level through a slider, giving them greater control over the system to tailor support according to students' programming abilities.

### 7.2 Limitations and Future Work

Our research has several limitations. First, our user study mainly focused on a fixed age group of 10- to 11-year-olds, and the children in our experiment already had a preliminary understanding of Scratch and basic programming concepts. While this experimental setup provided us with a controlled environment, eliminating potential cognitive differences among children of different ages, it may also limit the system's adaptability to children of other age groups, preventing us from identifying possible limitations. We encourage future studies to include children of various ages and interests, such as those aged 6 to 9.

Secondly, we compared the learning outcomes of MindScratch with traditional Scratch teaching. MindScratch leverages Generative AI technology to provide support at different stages of project development, thus improving learning outcomes. However, for children, the development of computational thinking is a long-term process. Currently, MindScratch does not offer assessments or feedback on completed programming projects, which limits its use as a self-learning tool. Future work could explore how to enhance MindScratch by integrating multi-perspective evaluation tools that

assess aesthetics, code quality, and computational thinking skills [4]. Additionally, MindScratch currently supports only Scratch, and it may have limitations when applied to more complex text-based programming languages like Python. Future studies could explore extending MindScratch to multiple programming languages to increase the system's applicability.

Finally, the safe use of AI by children has become a topic of public discussion. As we enter an era dominated by Generative AI, ethical considerations related to privacy and safety are increasingly important. For safety, we incorporated a moderation layer from OpenAI and prompt engineering to ensure that outputs contain no sexual content, hate speech, harassment, violence, or self-harm. However, this input-based filtering method may be influenced by biases in the training data. Future research could consider deploying large language models in specialized domains to limit the generation of uncontrollable responses, similar to our focus on programming support in this study. Such limitations could help reduce risks while still leveraging the capabilities of large language models.

## 8   Conclusion

This paper introduces MindScratch, a visual programming support tool that assists students in completing creative exploration and programming implementation in the classroom. We are the first to propose a GAI-driven interactive mind map system to help students achieve learning goals and offer timely creative and programming support. The novelty of our research centers on the design of a mind map, which is intelligent, interactive, and capable of reducing the burden of teachers by aligning learning goals, supporting students' programming processes, and providing multimodal resources. We compared MindScratch to Scratch through a within-subject study involving 24 participants. The results show that MindScratch facilitates participants in achieving learning goals, enhancing code quality and creativity, and improving their CT skills. Additionally, we interviewed six programming educators to gather insights into the future of AI-driven educational tools. Our work offers insights and design considerations for building programming support tools to assist children with creative programming learning.

## Acknowledgements

## Funding

## Disclosure statement

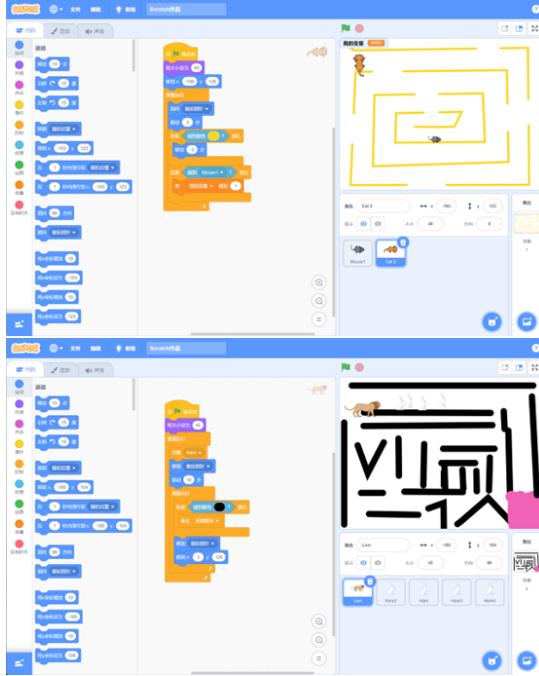No potential conflict of interest was reported by the author(s).

## References

[1] Teresa M Amabile. 1982. Social psychology of creativity: A consensual assessment technique. *Journal of personality and social psychology* 43, 5 (1982), 997.

[2] Brett A Becker, Paul Denny, James Finnie-Ansley, Andrew Luxton-Reilly, James Prather, and Eddie Antonio Santos. 2023. Programming is hard-or at least it used to be: Educational opportunities and challenges of ai code generation. In *Proceedings of the 54th ACM Technical Symposium on Computer Science Education V. 1*. 500–506.

[3] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems* 33 (2020), 1877–1901.

[4] Xiaolin Chai, Yan Sun, and Yan Gao. 2023. Towards Data-Driving Multi-View Evaluation Framework for Scratch. *Tsinghua Science and Technology* 29, 2 (2023), 517–528.

[5] Liuqing Chen, Shuhong Xiao, Yunnong Chen, Ruoyu Wu, Yaxuan Song, and Lingyun Sun. 2024. ChatScratch: An AI-Augmented System Toward Autonomous Visual Programming Learning for Children Aged 6-12. In *In: Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. 1–19. https://doi.org/10.1145/3613904.3642229

[6] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374* (2021).

[7] Erin Cherry and Celine Latulipe. 2014. Quantifying the creativity support of digital tools through the creativity support index. *ACM Transactions on Computer-Human Interaction (TOCHI)* 21, 4 (2014), 1–25.

[8] code.org 2023. code.org. https://code.org/ [Accessed on Dec. 27, 2023].

[9] Alina A Davier, Paul W Holland, and Dorothy T Thayer. 2004. *The kernel method of test equating*. Springer.

[10] Paul Denny, James Prather, Brett A Becker, James Finnie-Ansley, Arto Hellas, Juho Leinonen, Andrew Luxton-Reilly, Brent N Reeves, Eddie Antonio Santos, and Sami Sarsa. 2024. Computing education in the era of generative AI. *Commun. ACM* 67, 2 (2024), 56–67.

[11] Griffin Dietz, Jimmy K Le, Nadin Tamer, Jenny Han, Hyowon Gweon, Elizabeth L Murnane, and James A Landay. 2021. Storycoder: Teaching computational thinking concepts through storytelling in a voice-guided app for children. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–15.

[12] Griffin Dietz, Nadin Tamer, Carina Ly, Jimmy K Le, and James A Landay. 2023. Visual StoryCoder: A Multimodal Programming Environment for Children's Creation of Stories. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–16.

[13] James Finnie-Ansley, Paul Denny, Brett A Becker, Andrew Luxton-Reilly, and James Prather. 2022. The robots are coming: Exploring the implications of openai codex on introductory programming. In *Proceedings of the 24th Australasian Computing Education Conference*. 10–19.

[14] James Finnie-Ansley, Paul Denny, Andrew Luxton-Reilly, Eddie Antonio Santos, James Prather, and Brett A Becker. 2023. My ai wants to know if this will be on the exam: Testing openai's codex on cs2 programming exercises. In *Proceedings of the 25th Australasian Computing Education Conference*. 97–104.

[15] Emily Arteaga Garcia, João Felipe Pimentel, Zixuan Feng, Marco Gerosa, Igor Steinmacher, and Anita Sarma. 2023. How to Support ML End-User Programmers through a Conversational Agent. In *2024 IEEE/ACM 46th International Conference on Software Engineering (ICSE)*. IEEE Computer Society, 618–629.

[16] Deepanway Ghosal, Navonil Majumder, Ambuj Mehrish, and Soujanya Poria. 2023. Text-to-Audio Generation using Instruction Tuned LLM and Latent Diffusion Model. *arXiv preprint arXiv:2304.13731* (2023).

[17] Elena L Glassman, Lyla Fischer, Jeremy Scott, and Robert C Miller. 2015. Foobaz: Variable name feedback for student code at scale. In *Proceedings of the 28th Annual ACM Symposium on User Interface Software & Technology*. 609–617.

[18] Michal Gordon, Eileen Rivera, Edith Ackermann, and Cynthia Breazeal. 2015. Designing a relational social robot toolkit for preschool children to explore computational concepts. In *Proceedings of the 14th international conference on interaction design and children*. 355–358.

[19] hourofcode 2023. Hour of Code. https://hourofcode.com/ [Accessed on Dec. 23, 2023].

[20] Yu-Chang Hsu, Natalie Roote Irie, and Yu-Hui Ching. 2019. Computational thinking educational policy initiatives (CTEPI) across the globe. *TechTrends* 63 (2019), 260–270.

[21] Gwo-Jen Hwang, Shu-Yun Chien, and Wen-Shiang Li. 2021. A multidimensional repertory grid as a graphic organizer for implementing digital games to promote students' learning performances and behaviors. *British Journal of Educational Technology* 52, 2 (2021), 915–933.

[22] Theresia Devi Indriasari, Andrew Luxton-Reilly, and Paul Denny. 2020. A review of peer code review in higher education. *ACM Transactions on Computing Education (TOCE)* 20, 3 (2020), 1–25.

[23] Peiling Jiang, Jude Rayan, Steven P Dow, and Haijun Xia. 2023. Graphologue: Exploring large language model responses with interactive diagrams. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*. 1–20.

[24] Mukund Kapoor. 2023. Negative Prompts in Stable Diffusion: A Beginner's Guide. https://www.greataiprompts.com/imageprompt/what-is-negative-prompt-in-stable-diffusion/ Accessed: 2023-10-25.

[25] Enkelejda Kasneci, Kathrin Seßler, Stefan Küchemann, Maria Bannert, Daryna Dementieva, Frank Fischer, Urs Gasser, Georg Groh, Stephan Günnemann, Eyke Hüllermeier, et al. 2023. ChatGPT for good? On opportunities and challenges of large language models for education. *Learning and individual differences* 103 (2023), 102274.

[26] Majeed Kazemitabaar, Justin Chow, Carl Ka To Ma, Barbara J Ericson, David Weintrop, and Tovi Grossman. 2023. Studying the effect of ai code generators on supporting novice learners in introductory programming. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–23.

[27] Majeed Kazemitabaar, Xinying Hou, Austin Henley, Barbara Jane Ericson, David Weintrop, and Tovi Grossman. 2023. How novices use LLM-based code generators to solve CS1 coding tasks in a self-paced learning environment. In *Proceedings of the 23rd Koli Calling International Conference on Computing Education Research*. 1–12.

[28] Frank Klassner and Scott D Anderson. 2003. Lego MindStorms: Not just for K-12 anymore. *IEEE robotics & automation magazine* 10, 2 (2003), 12–18.

[29] Özgen Korkmaz and BAI Xuemei. 2019. Adapting computational thinking scale (CTS) for Chinese high school students and their thinking scale skills level. *Participatory Educational Research* 6, 1 (2019), 10–26.

[30] Anastasia Kovalkov, Benjamin Paaßen, Avi Segal, Niels Pinkwart, and Kobi Gal. 2021. Automatic creativity measurement in scratch programs across modalities. *IEEE Transactions on Learning Technologies* 14, 6 (2021), 740–753.

[31] Chang Liu, Loc Hoang, Andrew Stolman, and Bo Wu. 2024. HiTA: A RAG-Based Educational Platform that Centers Educators in the Instructional Loop. In *International Conference on Artificial Intelligence in Education*. Springer, 405–412.

[32] Kirsti Lonka, Juho Makkonen, Minna Berg, Markus Talvio, Erika Maksniemi, Milla Kruskopf, Heidi Lammassaari, Lauri Hietajärvi, and Suvi Krista Westling. 2018. *Phenomenal learning from Finland*. Edita, Finland.

[33] Stephen MacNeil, Andrew Tran, Dan Mogil, Seth Bernstein, Erin Ross, and Ziheng Huang. 2022. Generating diverse code explanations using the gpt-3 large language model. In *Proceedings of the 2022 ACM Conference on International Computing Education Research-Volume 2*. 37–39.

[34] Orni Meerbaum-Salant, Michal Armoni, and Mordechai Ben-Ari. 2011. Habits of Programming in Scratch. In *Proceedings of the 16th Annual Joint Conference on Innovation and Technology in Computer Science Education* (Darmstadt, Germany) *(ITiCSE '11)*. Association for Computing Machinery, New York, NY, USA, 168–172. https://doi.org/10.1145/1999747.1999796

[35] Orni Meerbaum-Salant, Michal Armoni, and Mordechai Ben-Ari. 2011. Habits of programming in scratch. In *Proceedings of the 16th annual joint conference on Innovation and technology in computer science education*. 168–172.

[36] Joseph Bahman Moghadam, Rohan Roy Choudhury, HeZheng Yin, and Armando Fox. 2015. AutoStyle: Toward coding style feedback at scale. In *Proceedings of the Second (2015) ACM Conference on Learning@ Scale*. 261–266.

[37] Jesús Moreno-León and Gregorio Robles. 2015. Dr. Scratch: A web tool to automatically evaluate Scratch projects. In *Proceedings of the workshop in primary and secondary computing education*. 132–133.

[38] Julie Mueller, Danielle Beckett, Eden Hennessey, and Hasan Shodiev. 2017. Assessing computational thinking across the curriculum. *Emerging research, practice, and policy on computational thinking* (2017), 251–267.

[39] Stable Diffusion Online. 2023. Stable Diffusion: A Latent Text-to-Image Diffusion Model. https://stablediffusionweb.com/ Accessed: 2023-09-08.

[40] Zhenhui Peng, Yuzhi Liu, Hanqi Zhou, Zuyu Xu, and Xiaojuan Ma. 2022. CReBot: Exploring interactive question prompts for critical paper reading. *International Journal of Human-Computer Studies* 167 (2022), 102898.

[41] Matei-Dan Popovici. 2023. ChatGPT in the classroom. Exploring its potential and limitations in a functional programming course. *International Journal of Human–Computer Interaction* (2023), 1–12.

[42] Dylan J Portelance and Marina Umaschi Bers. 2015. Code and Tell: Assessing young children's learning of computational thinking using peer video interviews with ScratchJr. In *Proceedings of the 14th international conference on interaction design and children*. 271–274.

[43] JD Zamfirescu-Pereira Laryn Qi, Bjoern Hartmann, and John DeNero Narges Norouzi. 2023. Conversational Programming with LLM-Powered Interactive Support in an Introductory Computer Science Course. (2023).

[44] Michael Reilly. 2013. The kindergarten coders. *New scientist* 2927 (2013), 21–22.

[45] Mitchel Resnick, John Maloney, Andrés Monroy-Hernández, Natalie Rusk, Evelyn Eastmond, Karen Brennan, Amon Millner, Eric Rosenbaum, Jay Silver, Brian Silverman, et al. 2009. Scratch: programming for all. *Commun. ACM* 52, 11 (2009), 60–67.

[46] Mitchel Resnick, John Maloney, Andrés Monroy-Hernández, Natalie Rusk, Evelyn Eastmond, Karen Brennan, Amon Millner, Eric Rosenbaum, Jay Silver, Brian Silverman, and Yasmin Kafai. 2009. Scratch: Programming for All. *Commun. ACM* 52, 11 (nov 2009), 60–67. https://doi.org/10.1145/1592761.1592779

[47] Andoni Rivera-Pinto, Johan Kildal, and Elena Lazkano. 2023. Toward programming a collaborative robot by interacting with its digital twin in a mixed reality environment. *International Journal of Human–Computer Interaction* (2023), 1–13.

[48] Anthony Robins, Janet Rountree, and Nathan Rountree. 2003. Learning and teaching programming: A review and discussion. *Computer science education* 13, 2 (2003), 137–172.

[49] Margarida Romero, Therese Laferriere, and Thomas Michael Power. 2016. The move is on! From the passive multimedia learner to the engaged co-creator. *ELearn* 2016, 3 (2016).

[50] Margarida Romero, Alexandre Lepage, and Benjamin Lille. 2017. Computational thinking development through creative programming in higher education. *International Journal of Educational Technology in Higher Education* 14 (2017), 1–15.

[51] Sherry Ruan, Liwei Jiang, Justin Xu, Bryce Joe-Kun Tham, Zhengneng Qiu, Yeshuang Zhu, Elizabeth L Murnane, Emma Brunskill, and James A Landay. 2019. Quizbot: A dialogue-based adaptive learning system for factual knowledge. In *Proceedings of the 2019 CHI conference on human factors in computing systems*. 1–13.

[52] Sami Sarsa, Paul Denny, Arto Hellas, and Juho Leinonen. 2022. Automatic generation of programming exercises and code explanations using large language models. In *Proceedings of the 2022 ACM Conference on International Computing Education Research-Volume 1*. 27–43.

[53] Cynthia Selby and John Woollard. 2013. Computational thinking: the developing definition. (2013).

[54] Karim Shabani, Mohamad Khatib, and Saman Ebadi. 2010. Vygotsky's zone of proximal development: Instructional implications and teachers' professional development. *English language teaching* 3, 4 (2010), 237–248.

[55] Murray Shanahan, Kyle McDonell, and Laria Reynolds. 2023. Role play with large language models. *Nature* 623, 7987 (2023), 493–498.

[56] Valerie J. Shute, Chen Sun, and Jodi Asbell-Clarke. 2017. Demystifying computational thinking. *Educational Research Review* 22 (2017), 142–158. https://doi.org/10.1016/j.edurev.2017.09.003

[57] Simon, Judy Sheard, Michael Morgan, Andrew Petersen, Amber Settle, Jane Sinclair, Gerry Cross, and Charles Riedesel. 2016. Negotiating the maze of academic integrity in computing education. In *Proceedings of the 2016 ITiCSE working group reports*. 57–80.

[58] Rishabh Singh, Sumit Gulwani, and Armando Solar-Lezama. 2013. Automated feedback generation for introductory programming assignments. In *Proceedings of the 34th ACM SIGPLAN conference on Programming language design and implementation*. 15–26.

[59] Harry Stokhof, Bregje De Vries, Theo Bastiaens, and Rob Martens. 2020. Using mind maps to make student questioning effective: Learning outcomes of a principle-based scenario for teacher guidance. *Research in Science Education* 50, 1 (2020), 203–225.

[60] Sangho Suh, Jian Zhao, and Edith Law. 2022. Codetoon: Story ideation, auto comic generation, and structure mapping for code-driven storytelling. In *Proceedings of the 35th Annual ACM Symposium on User Interface Software and Technology*. 1–16.

[61] Xiaohang Tang, Xi Chen, Sam Wong, and Yan Chen. 2023. VizPI: A Real-Time Visualization Tool for Enhancing Peer Instruction in Large-Scale Programming Lectures. In *Adjunct Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*. 1–3.

[62] Christina Tikva and Efthimios Tambouris. 2021. Mapping computational thinking through programming in K-12 education: A conceptual model based on a systematic literature Review. *Computers & Education* 162 (2021), 104083.

[63] Joke Voogt, Petra Fisser, Jon Good, Punya Mishra, and Aman Yadav. 2015. Computational thinking in compulsory education: Towards an agenda for research and practice. *Education and information technologies* 20 (2015), 715–728.

[64] April Yi Wang, Yan Chen, John Joon Young Chung, Christopher Brooks, and Steve Oney. 2021. Puzzleme: Leveraging peer assessment for in-class programming exercises. *Proceedings of the ACM on Human-Computer Interaction* 5, CSCW2 (2021), 1–24.

[65] April Yi Wang, Andrew Head, Ashley Ge Zhang, Steve Oney, and Christopher Brooks. 2023. Colaroid: A literate programming approach for authoring explorable multi-stage tutorials. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–22.

[66] Xinyuan Wang, Qian Xing, Qiao Jin, and Danli Wang. 2024. "be a lighting programmer": Supporting children collaborative learning through tangible programming system. *International Journal of Human–Computer Interaction* 40, 10 (2024), 2622–2640.

[67] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems* 35 (2022), 24824–24837.

[68] Tongshuang Wu, Michael Terry, and Carrie Jun Cai. 2022. Ai chains: Transparent and controllable human-ai interaction by chaining large language model prompts. In *Proceedings of the 2022 CHI conference on human factors in computing systems*. 1–22.

[69] Lixiang Yan, Lele Sha, Linxuan Zhao, Yuheng Li, Roberto Martinez-Maldonado, Guanliang Chen, Xinyu Li, Yueqiao Jin, and Dragan Gašević. 2024. Practical and ethical challenges of large language models in education: A systematic scoping

review. *British Journal of Educational Technology* 55, 1 (2024), 90–112.

[70] Tzu-Chi Yang and Zhi-Shen Lin. 2024. Enhancing elementary school students' computational thinking and programming learning with graphic organizers. *Computers & Education* 209 (2024), 104962.

[71] Jia-Yu Yao, Kun-Peng Ning, Zhen-Hui Liu, Mu-Nan Ning, and Li Yuan. 2023. Llm lies: Hallucinations are not bugs, but features as adversarial examples. *arXiv preprint arXiv:2310.01469* (2023).

[72] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. 2022. ReAct: Synergizing Reasoning and Acting in Language Models. In *The Eleventh International Conference on Learning Representations*.

[73] Kangyu Yuan, Hehai Lin, Shilei Cao, Zhenhui Peng, Qingyu Guo, and Xiaojuan Ma. 2023. CriTrainer: An Adaptive Training Tool for Critical Paper Reading. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*. 1–17.

[74] Ashley Ge Zhang, Yan Chen, and Steve Oney. 2023. Vizprog: Identifying misunderstandings by visualizing students' coding progress. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–16.

[75] Chao Zhang, Cheng Yao, Jiayi Wu, Weijia Lin, Lijuan Liu, Ge Yan, and Fangtian Ying. 2022. StoryDrawer: A Child–AI Collaborative Drawing System to Support Children's Creative Visual Storytelling. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*. 1–15.

[76] Lvmin Zhang and Maneesh Agrawala. 2023. Adding Conditional Control to Text-to-Image Diffusion Models. arXiv:2302.05543 [cs.CV]

[77] Li Zhao, Xiaohong Liu, Chenhui Wang, and Yu-Sheng Su. 2022. Effect of different mind mapping approaches on primary school students' computational thinking skills during visual programming learning. *Computers & Education* 181 (2022), 104445.

## A  Appendix A

**Table 7: Detailed descriptions of the two theme-based creative programming tasks.**

| N | Project Description | Learning Objectives | The Scratch Code and Creations |
|---|---|---|---|
| 1 | Cat and Mouse: Create a maze background with two characters, a cat and a mouse. When the cat touches the mouse, the cat grows larger. Both the cat and the mouse can eat cheese, and the mouse will randomly appear at any location in the maze. | (1) Students can proficiently use various methods to control the cat or mouse. (2) Students can use relevant code for detection and decision-making functions to complete the game's features. |  |
| 2 | Kitten Fishing: Create a game where a kitten is fishing in a pond. There are four characters: the kitten, the fishing rod, the fish, and the score. When the fishing rod catches a fish, the score increases by one. | (1) First, set up the pond stage. (2) Second, make the fish appear anywhere in the pond and enable them to move back and forth in the water. (3) Third, control the fishing rod in the kitten's hand to move downwards. If the rod touches a fish, the fish should disappear, and the score counter should increase by one. |  |