

Scaffolding Metacognition in Programming Education: Understanding Student–AI Interactions and Design Implications

BOXUAN MA*, Kyushu University, Japan

HUIYONG LI, Kyushu University, Japan

GEN LI, Kyushu University, Japan

LI CHEN, Osaka Kyoiku University, Japan

CHENG TANG, Kyushu University, Japan

YINJIE XIE, Kyushu University, Japan

CHENGHAO GU, Kyushu University, Japan

ATUSHI SHIMADA, Kyushu University, Japan

SHIN'ICHI KONOMI, Kyushu University, Japan

Generative AI tools such as ChatGPT now provide novice programmers with unprecedented access to instant, personalized support. While this holds clear promise, their influence on students' metacognitive processes remains underexplored. Existing work has largely focused on correctness and usability, with limited attention to whether and how students' use of AI assistants supports or bypasses key metacognitive processes. This study addresses that gap by analyzing student–AI interactions through a metacognitive lens in university-level programming courses. We examined more than 10,000 dialogue logs collected over three years, complemented by surveys of students and educators. Our analysis focused on how prompts and responses aligned with metacognitive phases and strategies. Synthesizing these findings across data sources, we distill design considerations for AI-powered coding assistants that aim to support rather than supplant metacognitive engagement. Our findings provide guidance for developing educational AI tools that strengthen students' learning processes in programming education.

CCS Concepts: • **Human-centered computing** → **Interactive systems and tools**; • **Social and professional topics** → **Computing education**.

Additional Key Words and Phrases: programming education, generative AI, metacognitive theory, design guidelines

ACM Reference Format:

Boxuan Ma, Huiyong Li, Gen Li, Li Chen, Cheng Tang, Yinjie Xie, Chenghao Gu, Atushi Shimada, and Shin'ichi Konomi. 2018. Scaffolding Metacognition in Programming Education: Understanding Student–AI Interactions and Design Implications. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, 29 pages. <https://doi.org/XXXXXXX.XXXXXXX>

*Corresponding author. Email: boxuan@artsci.kyushu-u.ac.jp

Authors' Contact Information: Boxuan Ma, boxuan@artsci.kyushu-u.ac.jp, Kyushu University, Fukuoka, Japan; Huiyong Li, Kyushu University, Fukuoka, Japan; Gen Li, Kyushu University, Fukuoka, Japan; Li Chen, Osaka Kyoiku University, Osaka, Japan; Cheng Tang, Kyushu University, Fukuoka, Japan; Yinjie Xie, Kyushu University, Fukuoka, Japan; Chenghao Gu, Kyushu University, Fukuoka, Japan; Atushi Shimada, Kyushu University, Fukuoka, Japan; Shin'ichi Konomi, Kyushu University, Fukuoka, Japan.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.

Manuscript submitted to ACM

Manuscript submitted to ACM

1

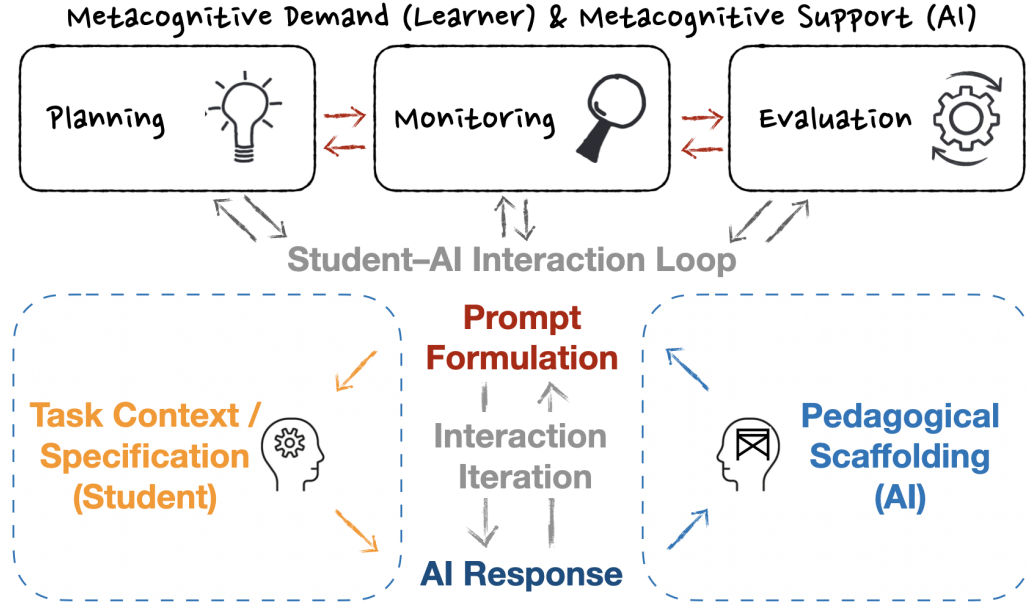


Fig. 1. Metacognitive demand (learner) and support (AI) in the student-AI interaction loop.

1 Introduction

Students learning to program routinely encounter deep uncertainty: they may stare at a blank editor unsure how to begin, struggle to decipher cryptic error messages when code fails, and—even when it runs—doubt whether their solution is efficient, generalizable, or aligned with best practices. These moments often arise precisely when help is hardest to access: office hours may not align with schedules, online forums can feel intimidating, and in crowded classrooms attention tends to go to the most vocal students, leaving feedback uneven and delayed [54, 55]. In higher-education programming courses, large enrollments, wide variation in prior experience, and limited instructional resources further constrain the timely, individualized support novices need [54, 55]. Such delays or generic feedback are linked to shallow understanding, weaker retention, limited transfer, and underdeveloped higher-order learning skills [39]. Compounding these constraints, many beginners lack a firm grasp of fundamentals (e.g., syntax, control flow, algorithms) and the metacognitive strategies required for planning, monitoring, and evaluation. Addressing this gap calls for scalable, equitable, and learner-centered approaches that deliver guidance at the moment students encounter difficulty.

Large language models (LLMs) and AI-powered coding assistants offer a promising avenue for addressing this need. They can instantly generate explanations, examples, and guidance, providing help precisely when students encounter difficulties [58, 61]. Recent work highlights how LLMs can augment students' workflows and make assistance broadly available outside lecture time or office hours [28, 56]. These capabilities make AI a compelling tool for addressing the long-standing gap between when students encounter difficulty and when they receive help. At the same time, this strength of AI tools introduces new pedagogical risks. Because AI can often produce working solutions with minimal prompting [12, 17], students may offload problem-solving steps to the AI assistants rather than engage in metacognitive processes such as planning, monitoring, and evaluation [2, 47]. Scholars have warned such “metacognitive laziness”

when learners rely on AI to do the thinking for them [11, 16]. These concerns extend to academic integrity and to the longer-term development of independent problem solving and metacognitive growth [53, 59].

In response, researchers and designers have begun to build constrained, pedagogically minded assistants that avoid simply handing back finished code. Systems such as CodeHelp [32] and CodeAid [28] exemplify interventions that restrict direct answers and instead provide pseudocode, hints, or explanation-first interactions to encourage student reasoning. These guardrails position AI as a partner in thinking rather than a replacement. However, prior work has largely emphasized use patterns and performance outcomes (e.g., correctness, time to solution), with comparatively little attention to metacognitive support. Without designs that align usage with metacognitive phases, such tools risk falling short of their educational potential. We argue that metacognition, understood as the ability to monitor and regulate one's own thinking, offers a productive lens for improving AI-based assistance. As Figure 1 shows, each turn in the student-AI interaction loop places metacognitive demands on learners and is best met with metacognitive support from the AI. Framing the exchange as a cycle that moves from planning to monitoring to evaluation clarifies when learners should plan, how they should monitor, and how they should evaluate, enabling assistants that strengthen rather than bypass metacognitive skills.

This paper addresses this gap by bringing metacognitive theory into conversation with empirical studies of real-world student-AI interactions in programming education. We present a multi-method study of university-level programming courses that combines analysis of 10,632 student-AI dialogue logs collected over three years with student surveys and educator interviews. This setup lets us observe how students naturally use AI assistance (without externally imposed constraints), how those interaction patterns align with metacognitive phases (planning, monitoring, evaluation), and what design principles derive for learner-centered AI scaffoldings that foster metacognitive skill development.

Guided by this motivation, our study addresses the following research questions:

- **RQ1 - Metacognitive Interaction Patterns:** What metacognitive engagement patterns (e.g., planning, monitoring, evaluation) do students exhibit when they interact with AI-powered assistants during programming tasks, and what response types do the assistants provide?
- **RQ2 - Student Perceptions:** How do students perceive AI support for programming learning, including perceived benefits, concerns, and suggestions for improvement?
- **RQ3 - Educator's Views:** How do educators view students' metacognitive prompts and the pedagogical potential (and risks) of AI-powered scaffolding in programming courses?
- **RQ4 - Design Principles:** What design principles for pedagogical AI assistants follow from the identified interaction patterns and stakeholder perceptions?

Synthesizing answers to our research questions, this paper offers a critical account of the design space for pedagogical AI assistants in programming education. Our analysis illuminates how students interact with AI-powered tools, revealing connections to metacognitive practices. From these insights, we derive actionable design principles for pedagogical AI systems that move beyond solution provision toward fostering planning, monitoring, and reflective evaluation rather than enabling metacognitive offloading.

2 Related Work

2.1 Capabilities of LLMs in Programming Education

LLMs are reshaping programming education by lowering barriers to help-seeking and automating a broad range of programming tasks, including code generation, debugging, code explanation, and the delivery of personalized feedback [12, 47].

Early research in this area has primarily focused on evaluating the capabilities of LLMs on programming tasks [18, 44, 47, 49, 50], showing that they are able to address common programming challenges with impressive results. For instance, evaluations of OpenAI Codex and GPT models showed that they could solve many standard CS1 and CS2 exam problems, even outperforming median student performance [17]. Extending this line of work, Savelka et al. [50] evaluated GPT-3 and GPT-4 on 599 programming exercises from three Python courses and found that while GPT-3 failed most assessments, GPT-4 was able to pass entire courses without human assistance. Similarly, Phung et al. [44] compared GPT models with human tutors and found that they achieved near human-level performance on Python programming tasks as well as in repairing real-world buggy programs. Other studies highlights LLMs' ability to generate not only functional code but also clear, structured explanations that improve learners' comprehension. For instance, MacNeil et al. [37] found that GPT-3's explanations helped students reason through problems more effectively. Similarly, Leinonen et al. [30] showed that LLM-generated explanations were often clearer and more pedagogically useful than student-written ones, though they occasionally suffered from inaccuracy or overconfidence. In the context of feedback, LLMs have shown promise in supporting students by commenting on programming assignments and facilitating the application of theoretical knowledge in practice. Prior work demonstrated that students generally evaluated personalized AI-generated feedback positively [41]. Zhang et al. [62] reported that learners in a CS1 Java course generally perceived AI-generated feedback as aligned with the established principles of formative feedback.

2.2 Student and Educator Perspectives with LLMs

While early research on LLMs in programming education has primarily emphasized model capabilities, a growing body of work has shifted attention toward the human side. Recent studies have begun to systematically examine how students and educators perceive and experience the use of LLMs in programming education [5, 21, 35, 36, 51, 61]. This is critical because the effectiveness of LLMs as educational tools depends not only on what they can do, but also on how learners and teachers engage with them and adapt their practices accordingly.

Studies consistently show that students view LLMs as valuable supports in programming education [35, 51, 61]. They appreciate the immediacy and accessibility of LLMs, emphasizing how quickly it can provide solutions or guidance compared to searching online resources or waiting for human assistance [5, 61]. Beyond efficiency, students also highlight the clarity of LLM-generated feedback, which helps them understand concepts and debug code more effectively [35]. At the same time, students remain cautious: many acknowledge that answers can be incorrect or misleading [35, 51]. They also worried about the unintended consequences of unmoderated LLM use lead to superficial understanding and lower knowledge retention [53, 59]. In the absence of deliberate practice, these behaviors risk undermining the development of problem-solving persistence, debugging proficiency, and critical thinking [25, 41, 51].

Educators' perspectives on LLMs have also been widely examined in prior studies. For instructors, such tools can reduce repetitive support work and free time for higher-level teaching, while also providing a private channel for students reluctant to ask questions in public [19]. Although LLMs are seen as offering unprecedented scalability and personalization, educators tend to approach them with caution [14, 25]. Compared to students, educators are

generally more skeptical, voicing stronger concerns about over-reliance, misuse, and ethical implications, and calling for institutional policies and guidelines [7]. Empirical studies further indicate that teaching practices must adapt as LLMs enter classrooms [1], and highlight the need for educators to remain abreast of technological developments while guiding students on ethical use [46]. Some advocate restricting or banning LLM use to preserve academic integrity [29], whereas others experiment with integration strategies that leverage AI for formative support while safeguarding assessment integrity [3]. Across these perspectives, there is broad agreement on the importance of pedagogical scaffoldings that encourage strategic help-seeking and foster metacognitive engagement [14, 28, 56].

2.3 LLM-based Coding Assistants

Concerns about student over-reliance on LLMs, together with evidence that many students struggle to write effective prompts and thus receive unhelpful feedback [31, 36], have led researchers to create custom tools that balance LLM capabilities with pedagogical goals. These tools often incorporate mechanisms such as stepwise hints [26], prompt construction for code generation [13], or restrictions on direct answers [28, 32], ensuring that AI assistance supports learning rather than bypassing it.

For example, Kazemitabaar et al. [26] created Coding Steps, a tool where students give prompts and receive code automatically generated by LLMs. Yet, they observed that many students simply copied the exercise text as prompts and accepted the AI output without modification. In response to such tendencies, Lifton et al. [32] developed CodeHelp, which supports programming learners by offering scaffolded guidance rather than full solutions. Students can submit questions alongside their code and optional error messages, and the system provides structured hints instead of direct answers. Building further on this idea, CodeAid [28] incorporates pre-defined templates, enabling features such as conceptual explanations, pseudo-code generation with step-by-step commentary, and annotated feedback on incorrect code.

In sum, current research has made strong progress in building guardrails that keep AI from generating answers directly. However, most studies stop at measuring performance or usability, leaving open questions about how these tools actually shape students' high-order thinking. In particular, there is little evidence on whether current designs truly support metacognitive development including planning, monitoring, and evaluation. This gap points to the need for research that examines how student interactions with AI align with and support metacognitive processes.

2.4 Metacognition and Self-regulated Learning in Programming Learning

Metacognitive skills have long been recognized as essential ability for programming learners [34], where novices often struggle with structuring approaches (planning) [15, 20], identifying and fixing errors (monitoring) [42], and optimizing solutions (evaluation) [48]. Prior research has consistently demonstrated the benefits of metacognitive support for learning engagement and outcomes. For example, Choi et al. [10] found that prompting reflection after programming tasks improved performance in both immediate and delayed tests. Similarly, Yilmaz and Yilmaz [23] reported higher engagement when students received personalized metacognitive feedback, and Cheng et al. [9] showed that high-performing students relied on elaboration and critical thinking, while lower-performing peers used more basic strategies.

In the context of human-AI interaction, however, concerns have emerged that easy access to AI-generated solutions may reduce students' engagement in metacognitive practices, a phenomenon described as "metacognitive laziness" [16]. Studies suggest that students sometimes bypass essential metacognitive processes such as evaluation when using AI-powered tools [8]. While scaffolds such as step-by-step pseudo-code can encourage reflection [28], overly complete

solutions may discourage persistence and exploration [26]. This tension highlights a critical gap: although the benefits of the metacognition for programming learning are well established, far less is known about how the metacognition phases and strategies evolve when learners interact with AI-powered assistants [60]. Since metacognition is a dynamic process, strategies may shift in response to AI-mediated interaction. Understanding these dynamics is essential for designing AI-powered pedagogical assistants that reinforce rather than replace students’ metacognitive engagement [43]. This gap motivates our study, which analyzes student–AI dialogues to uncover key metacognitive strategies and derive actionable design principles for pedagogical AI systems.

3 Method

We conducted a multi-year study in an introductory undergraduate Python programming course at a national university in Japan. The study spanned from 2023 to 2025 and consisted of four offerings of the same course (one offering in 2023, two in 2024, and one in 2025), with a total of 248 students. Throughout each course, all students were given optional access to an AI tool as a supplementary learning resource. Access was provided through a custom, web-based interface that enabled real-time interaction with GPT-based AI. All conversations between students and AI were systematically logged for subsequent analysis. The study adopted a mixed-methods approach, combining quantitative and qualitative data collection. The study received approval from the university’s ethics review board prior to the start of the study.

3.1 Course Structure

This first-year undergraduate course, designed for beginners, covered the foundations of the Python programming language in 14 lessons delivered over one semester. Each 90-minute lesson was divided into two segments: the first 45 minutes were dedicated to direct instruction, supported by lecture presentations and slides, followed by 45 minutes of hands-on programming tasks related to the lesson’s topic (e.g., write a specific function based on the requirements). This structure was intended to enable students to immediately apply the theoretical concepts introduced in class.

Two instructors delivered the four course offerings, and the curriculum was largely unchanged between 2023 and 2025, which allowed us to maintain consistency across offerings. In each lesson, students completed programming exercises. Each course included a total of 57 exercises (averaging 4.75 per lesson). All exercises were submitted through the university’s Learning Management System (LMS), which automatically logged and evaluated each exercise submission, generating a complete and traceable record of students’ programming activity. Students had access to a variety of learning resources beyond AI, including lecture slides and support from the course instructors, and were instructed to cite any external sources used in their work.

3.2 Participants

The study involved a total of 248 undergraduate students over a three-year period (62 participants in 2023, 126 in 2024, and 60 in 2025), across four offerings of the same course in the university. Participation was voluntary, informed consent was obtained from all participants, and no aspect of participation affected students’ course grades. The distribution of demographic characteristics was generally consistent across years. Overall, the majority of participants (96.8%) were first-year students at the time of participation, with a near-equal gender distribution (53.2% male, 45.8% female, and 1% reporting a gender identity outside the male–female binary). While most participants were domestic students from Japan, approximately 14.5% were international students from diverse countries.

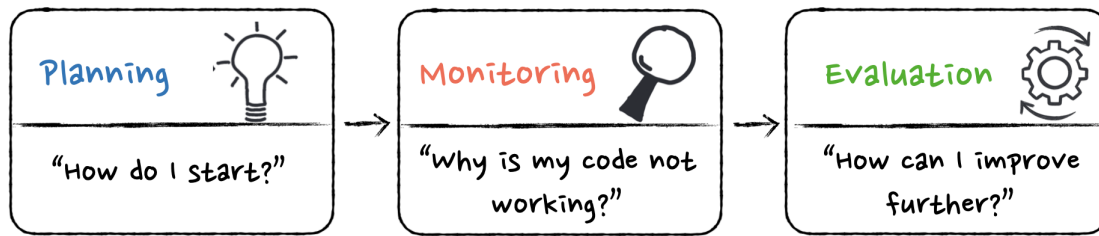


Fig. 2. Metacognitive Phases of Student–AI Interaction in Programming Tasks

3.3 Data Sources

To gather a comprehensive understanding of student interactions with AI and student and educator perceptions of AI, we employed a multifaceted data collection approach by collecting student-AI interaction logs, pre-questionnaire and post-questionnaire from students.

A primary data source for understanding students’ usage patterns and AI’s response was student-AI dialogue logs. Dialogue data has become an important data source to understand programming experiences [6] and coding approaches [18, 22], particularly when interacting with AI [26, 27]. For each question asked by students, we closely examined its content and AI generated responses through a thematic analysis.

At the beginning of the course, we asked all participants to fill out a pre-questionnaire to collect demographics such as major and gender. We also asked about their programming experiences and their familiarity with GenAI. As the course progressed, the students’ perceptions of using AI in programming education became more evident. To capture these insights, we conducted a questionnaire after the course, gathering students’ opinions to understand the advantages and challenges of integrating AI into the programming curriculum from their perspective. We also asked open questions about students’ viewpoints on the use of AI for programming learning purposes, how useful they found it, what they liked or disliked, and any open-ended feedback about it. All the questionnaire items are provided in Appendix B.

3.4 Thematic Analysis

Our analysis sought to capture the student-AI interactions reflect or support students’ metacognitive processes (e.g., planning, monitoring, evaluating). To this end, we conducted a thematic analysis of the dialogue logs. From the 10,632 recorded interactions, we performed stratified random sampling across cohorts, drawing 2,782 instances (26%) in total. Each year contributed a proportional segment of the sample, ensuring balanced representation for subsequent thematic analysis.

Drawing on coding schemes from prior studies [28, 36], we organized our analysis around two high-level dimensions: student prompts and AI-generated responses (see Table 1). This two-part structure enabled us to systematically capture both the student-facing and system-facing aspects of the interaction.

For the first dimension, our aim was to characterize not only the types of questions students asked but also the metacognitive phases embedded in these questions. While earlier codebooks had provided useful structures for describing metacognitive processes in programming contexts [45, 52], they had not been applied to the setting of student-AI interactions. To address this gap, we extended these prior frameworks to design a codebook that explicitly mapped prompts onto planning, monitoring, and evaluation activities, alongside related help-seeking categories. As shown in Figure 2, in the planning phase, students seek support in designing initial solution strategies, which includes help with

Table 1. The sub-dimensions from our thematic analysis, and inter-rater reliability metrics using Cohen’s Kappa and percentage agreement. The detailed codebook is provided in Appendix A.

Sub-Dimensions	Codes	Inter-Rater Reliability
What do students ask AI in <i>Planning</i> phase?	<i>Conceptual Question (CQ), Problem Understanding (PU), Asking for Example (AE), Code Implementation Questions (CI), Code Generation (CG)</i>	88% ($\kappa = .79$)
What do students ask AI in <i>Monitoring</i> phase?	<i>Error Message Interpretation (EM), Code Correction (CC), Code Verification (CV)</i>	90% ($\kappa = .85$)
What do students ask AI in <i>Evaluation</i> phase?	<i>Code Explanation (CE), Code Optimization (CO)</i>	100% ($\kappa = 1.0$)
How much is AI directly revealing the solution?	<i>Exact Solution Code, Steps to Fix Semantic Issue, Steps to Fix Syntax Issue, Step to Fix External Issue, Example Code, Conceptual Explanation, Code Explanation</i>	91% ($\kappa = .87$)
How technically correct is the response?	<i>Correct, Incorrect</i>	99% ($\kappa = .85$)
How helpful is the response if correct?	<i>Helpful, Not helpful</i>	100% ($\kappa = 1$)

understanding the problem requirements, exploring potential approaches, and generating starter code or conceptual examples that guide their entry into the task. During monitoring, students seek support in identifying and resolving issues during code execution, which involves assistance with interpreting error messages, verifying program correctness, and debugging both syntactic and semantic problems as they occur. In the evaluation phase, students seek support in reflecting on and improving code quality, which entails guidance on code explanation, optimization, and refinement, enabling learners to go beyond minimal functionality and enhance readability, efficiency, and overall robustness. This codebook allowed us to examine how students’ interactions with AI reflected underlying metacognitive engagement.

The second dimension examined AI-generated responses, emphasizing pedagogical quality and relevance. Building on categories adapted from prior work [28], we refined and consolidated a codebook and coded each reply on three independent axes: (1) Solution Revelation (ranging from exact solution code; steps to fix semantic, syntax, or external issues; to example code, conceptual explanation, and code explanations), (2) Technical Correctness (correct vs. incorrect), and (3) Helpfulness (helpful vs. not helpful). A key refinement was to treat out-of-scope solutions as Not Helpful even when technically correct—i.e., answers that recommend approaches beyond course expectations or constraints. These clarifications tightened category boundaries and aligned the scheme with our study context.

The coding procedure was identical across both dimensions. Following an inductive approach [4], two researchers jointly reviewed an initial set of 800 randomly sampled interactions, generating preliminary sub-dimensions and codes. They then independently coded an additional 150 samples using this preliminary codebook. Then, the results were discussed with the course instructors, and disagreements were resolved through consensus, leading to a refined version of the codebook. To establish reliability, the two researchers independently coded 334 further samples, and inter-rater agreement was calculated using Cohen’s Kappa and percentage agreement [38, 40]. Once reliability thresholds were satisfied, the researchers proceeded to independently code an additional 2,448 interactions randomly drawn from the remaining dataset. The finalized codebook is provided in Appendix A.

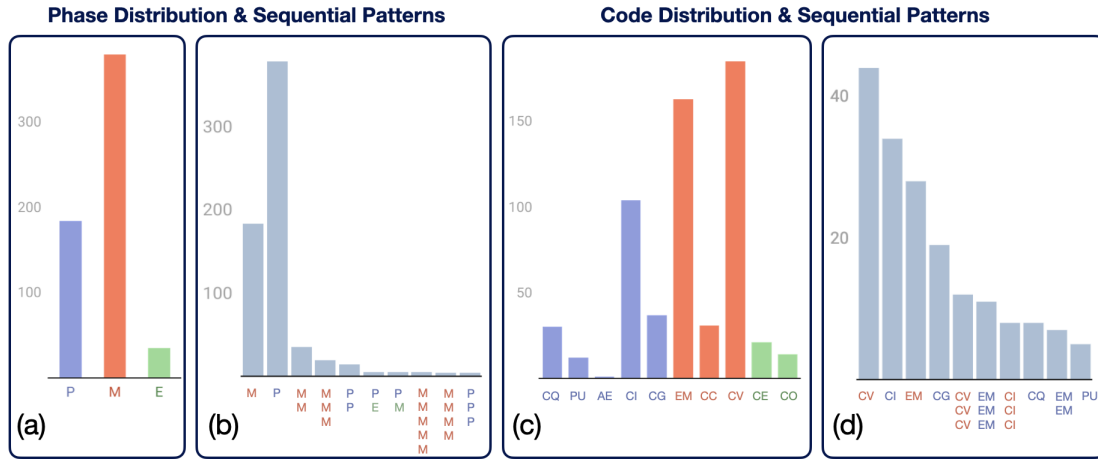


Fig. 3. Distribution and sequential patterns of student prompts across metacognitive phases and specific categories: (a) distribution of phases, (b) distribution of strategies, (c) sequential phase patterns, and (d) sequential strategy patterns.

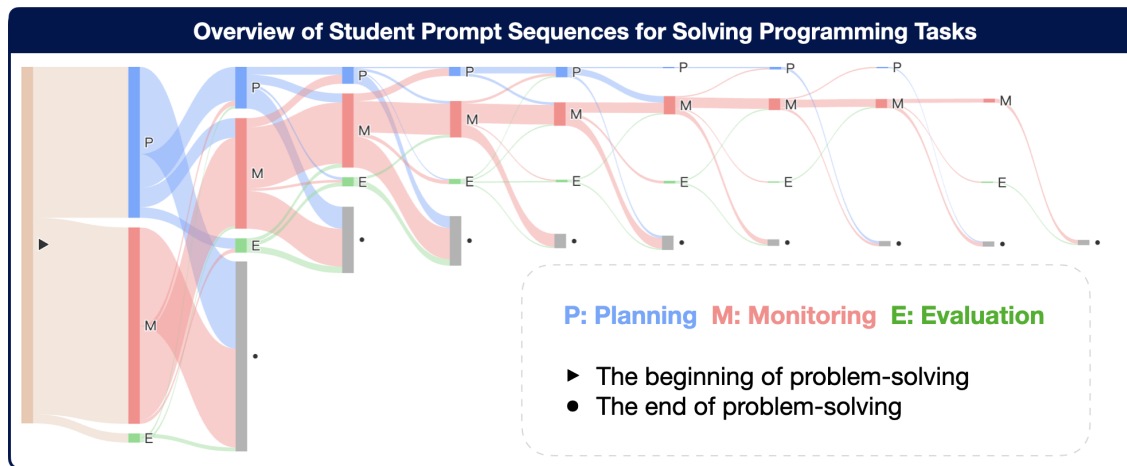


Fig. 4. Sequences of prompt during problem solving. Each path shows a per-problem trajectory from Start to End. Columns are successive turns. Nodes are grouped as P (Planning), M (Monitoring), and E (Evaluation).

4 Analysis of Interactions (RQ1)

4.1 Analysis of Student Prompts

To understand how students engaged with AI across different stages of problem solving, we first examined the distribution of prompts by metacognitive phase and strategies. Figure 3(a) presents the overall distribution of student prompts across the three metacognitive phases. Monitoring (M) accounted for the largest share of prompts, far exceeding Planning (P) and Evaluation (E). This indicates that students most frequently sought AI assistance during active debugging, troubleshooting, and correctness verification, rather than during initial strategy formulation or reflective evaluation. As shown in Figure 3(b), within the strategies, Error Message Interpretation (EM), Code Correction (CC), and Code

Verification (CV) dominated, highlighting students' tendency to rely on AI as a just-in-time debugging tool. In contrast, prompts associated with planning activities, such as Conceptual Questions (CQ), Problem Understanding (PU), or Asking for Example (AE), were comparatively sparse, and evaluation-related prompts (Code Explanation and Code Optimization) were least frequent. These descriptive patterns suggest that students primarily positioned AI as a reactive aid to overcome immediate coding obstacles, rather than as a proactive partner for strategic planning or reflective learning.

To examine how prompts unfolded over time, Figure 4 visualizes sequential flows with a Sankey diagram. A substantial share of sessions began directly in Monitoring rather than Planning, indicating that students often approached AI reactively after encountering execution problems rather than proactively during problem framing. Evaluation prompts often appeared as terminal moves, occurring at the end of problem-solving sequences, consistent with the notion that reflection was an afterthought rather than an integrated process. These patterns demonstrate that student-AI interactions tended to be short-sighted and narrowly focused on immediate code repair.

Figure 3 (c-d) further illustrates sequential patterns of student prompts in phases and strategies. At the phase level (Figure 3(c)), most frequent prompt sequences consist of a single phase, with monitoring are the most, followed by planning, while evaluation were rare. The majority of sequences were short and dominated by repetitive Monitoring, confirming the persistence of debugging-focused interactions. Longer or cross-phase sequences were comparatively rare, with only a few Planning-Monitoring-Evaluation trajectories, indicating that students seldom cycled through a full metacognitive loop with AI. At the strategy level (Figure 3(d)), the most frequent transitions involved Code Verification (CV), Code Implementation Question (CI), and Error Message Interpretation (EM), often chained in quick succession. These transitions suggest that students tended to oscillate between checking correctness, requesting implementation, and interpreting errors, without extending into broader conceptual or evaluative strategies. Collectively, the sequential analyses reinforce the picture of students treating AI primarily as an immediate error-fixing assistant, with limited movement toward higher-order reflection or iterative planning improvement.

Finally, Figure 5 quantifies these patterns using a Markov chain model. Monitoring showed the highest self-transition probability (0.49), confirming its stickiness. Planning also exhibited a self-loop (0.22), though weaker, suggesting some iteration in problem framing. Most cross-phase transitions converged on Monitoring, the dominant cross-phase pathway was Evaluation-to-Monitoring (0.54) and Planning-to-Monitoring (0.17). By contrast, transitions into Planning and Evaluation were weak, and nearly half of the sessions terminated after Planning (0.34) and Evaluation (0.46). Together, these results underscore that students positioned AI primarily as a reactive debugging assistant, with Planning under-specified and Evaluation marginal, leaving metacognitive phases fragmented rather than integrated.

4.2 AI-generated Responses

4.2.1 Response Types across Phases and Strategies. Table 2 details the types of responses the AI tended to produce across phases. In Planning, outputs were dominated by Example Code (S5, 48.3%) and Exact Solution Code (S1, 35.2%), with smaller proportions of Conceptual Explanation (S6, 14.2%) and very little Code Explanation (S7, 2.3%). This pattern aligns with the strategy-level distribution: most Planning prompts elicited examples or conceptual scaffolds, yet Code Generation (CG) consistently yielded runnable code, while Code Implementation (CI) was split—some responses provided direct solutions, but a larger share consisted of illustrative examples.

Monitoring exhibited a markedly different profile in which the model tended to “patch first.” Here, Exact Solution Code (S1) accounted for 51.7% of responses and Steps to Fix Syntax Issues (S3) for 30.3%. Within Monitoring categories,

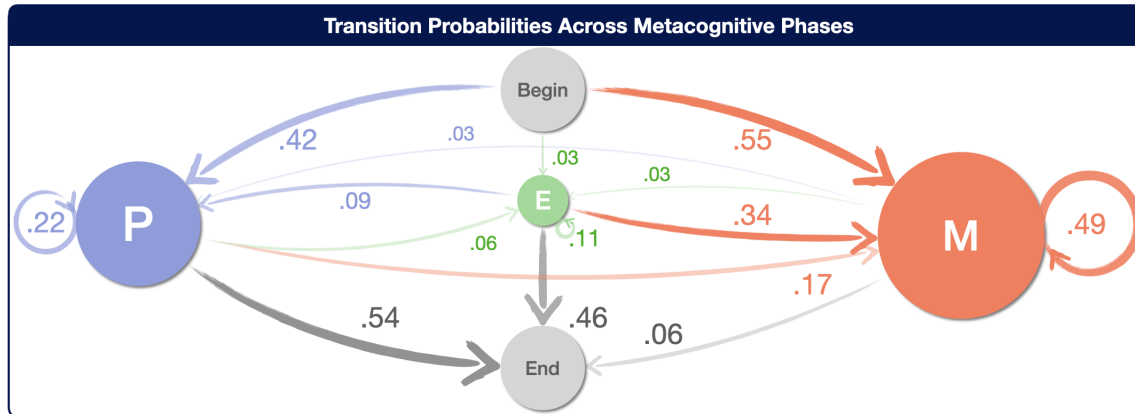


Fig. 5. Markov Chain Analysis of Phase Transitions.

Table 2. Distribution of response types (S1–7) across phases and strategies.

Phase	Category	S1	S2	S3	S4	S5	S6	S7
Planning (P)	Conceptual Question (CQ)	6.9%	0.0%	0.0%	0.0%	34.5%	51.7%	6.9%
	Problem Understanding (PU)	0.0%	0.0%	0.0%	0.0%	33.3%	50.0%	16.7%
	Asking for Example (AE)	0.0%	0.0%	0.0%	0.0%	100.0%	0.0%	0.0%
	Code Implementation Question (CI)	23.7%	0.0%	0.0%	0.0%	72.2%	4.1%	0.0%
	Code Generation (CG)	100.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%
	Overall	35.2%	0.0%	0.0%	0.0%	48.3%	14.2%	2.3%
Monitoring (M)	Error Message Interpretation (EM)	30.1%	5.1%	51.3%	1.9%	10.9%	0.6%	0.0%
	Code Correction (CC)	80.0%	3.3%	10.0%	0.0%	0.0%	0.0%	6.7%
	Code Verification (CV)	66.5%	0.0%	14.7%	0.0%	0.0%	1.2%	17.6%
	Overall	51.7%	2.5%	30.3%	0.8%	4.8%	0.8%	9.0%
Evaluation (E)	Code Explanation (CE)	0.0%	0.0%	5.3%	0.0%	15.8%	15.8%	63.2%
	Code Optimization (CO)	75.0%	0.0%	0.0%	0.0%	16.7%	8.3%	0.0%
	Overall	29.0%	0.0%	3.2%	0.0%	16.1%	12.9%	38.7%

Note. Response types: S1 = Exact Solution Code, S2 = Step to Fix Semantic Issue, S3 = Step to Fix Syntax Issue, S4 = Step to Fix External Issue, S5 = Example Code, S6 = Conceptual Explanation, S7 = Code Explanation.

Error Message Interpretation (EM) most often led to syntax-fix steps, whereas both Code Correction (CC) and Code Verification (CV) were dominated by direct runnable code.

In Evaluation, responses shifted toward interpretive support. The phase showed the highest share of Code Explanation (S7, 38.7%), alongside Exact Solution Code (S1, 29.0%), Example Code (S5, 16.1%), and Conceptual Explanation (S6, 12.9%). At the strategy-level, Code Explanation (CE) was almost entirely composed of Code Explanation (S7), while Code Optimization (CO) predominantly produced Exact Solution Code (S1).

4.2.2 Correctness of AI Responses. As shown in Table 3, the overall correctness rate of AI responses was high, but uneven across phases and categories. Planning responses (98.9%) were almost always correct, likely because conceptual clarifications and small examples leave less room for factual error. Evaluation responses (100%) similarly achieved perfect correctness, as explanation and commentary tasks required little speculative inference. By contrast, Monitoring

Table 3. Correctness and helpfulness rates by phase and category.

Phase	Category	Correct (%)	Helpful (%)
Planning (P)	Conceptual Question (CQ)	100.0%	89.7%
	Problem Understanding (PU)	100.0%	100.0%
	Asking for Example (AE)	100.0%	100.0%
	Code Implementation Question (CI)	97.9%	89.7%
	Code Generation (CG)	100.0%	100.0%
	Phase Overall	98.9%	92.6%
Monitoring (M)	Error Message Interpretation (EM)	95.5%	91.0%
	Code Correction (CC)	80.0%	80.0%
	Code Verification (CV)	95.3%	89.4%
	Phase Overall	94.1%	89.3%
Evaluation (E)	Code Explanation (CE)	100.0%	94.7%
	Code Optimization (CO)	100.0%	91.7%
	Phase Overall	100.0%	93.5%

responses (94.1%) displayed notable vulnerabilities, particularly in Code Correction (80% correct). A recurrent issue was overcorrection: the AI attempted to “improve” code by introducing extra modifications that were unnecessary or even harmful, sometimes converting a correct snippet into a faulty one. This failure mode underscores a fundamental challenge of LLMs in educational contexts—the tendency to produce overconfident edits when the minimal fix would suffice. The result is that while Monitoring responses were often technically competent, their reliability was uneven, making them less trustworthy for novice learners who lack the expertise to discriminate between accurate and overextended corrections.

4.2.3 Helpfulness of AI Responses. Helpfulness ratings (Table 3) revealed a second dimension beyond technical correctness: whether responses aligned with students’ learning needs. Planning responses (92.6%) and Evaluation responses (93.5%) were rated most helpful, since they provided clarifications, conceptual framing, and interpretive insights that students could directly use to advance their reasoning. Monitoring responses (89.3%), however, lagged behind. Even when technically correct, many were marked unhelpful when they supplied solutions beyond course scope (e.g., unnecessarily advanced algorithms or libraries) or when they required contextual details that students had not provided (e.g., incomplete I/O specifications). These cases illustrate a crucial distinction: correctness alone does not guarantee pedagogical usefulness. An answer can be “right” yet still fail to move the learner forward in a productive way. Conversely, answers judged unhelpful often reinforced passivity: by handing students runnable code, the AI risked bypassing the reasoning steps that are central to programming education. This reveals a pedagogical tension: effective AI support must balance accuracy with calibration to learner context, encouraging incremental progress without substituting for critical thinking.

5 Student Perspectives (RQ2)

5.1 Analysis of Pre-Questionnaire

The pre-questionnaire revealed that most participants (90%) self-identified as beginners in Python programming, and this proportion remained relatively consistent across all three years of the study. Reported usage of AI increased sharply

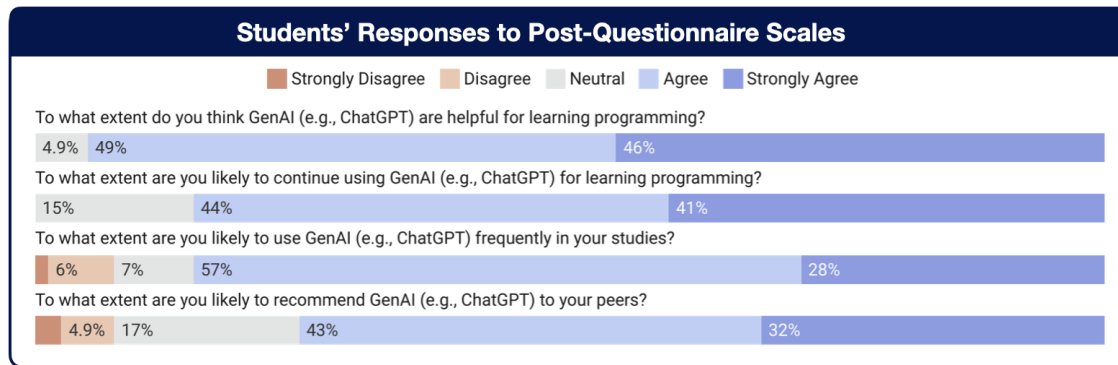


Fig. 6. Distribution of students' responses across all post-questionnaire scales.

over the three years of our study, reflecting both the rapid advancement of AI technology and its growing integration into students' everyday learning practices. In 2023, only 36.8% of respondents reported active use, while 57.9% indicated non-use. By 2024, active use had risen to 63.9%, with non-use dropping to 11.1%. By 2025, this trend became even more pronounced: 91.7% of respondents reported active use, and only 8.3% reported non-use. Similarly, students' self-perceived understanding of AI increased consistently across the study period. In 2023, only 5.2% of respondents indicated that they could explain AI in detail, 68.4% reported somewhat understanding, and 26.3% indicated little or no understanding. By 2024, 47.2% reported being able to explain AI in detail or to some extent, another 47.2% reported somewhat understanding, and only 5.6% indicated low understanding, with no students reporting no knowledge. In 2025, 100% of respondents indicated at least some understanding, and none reported low or no understanding. Finally, students' perceptions of the impact of AI on their learning were generally positive and increased across years. Overall, 82.3% of respondents considered AI to be beneficial or somewhat beneficial, with the proportion rising from 78.9% in 2023 to 80.5% in 2024, and 100% in 2025.

5.2 Analysis of Post-Questionnaire

Figure 6 presents an overview of student responses to the post-questionnaire. The results indicate strong support for AI as helpful tool in learning programming, with 45.7% completely agreeing and 49.4% agreeing. Similarly, the intention to continue using ChatGPT for learning programming remains high, with 40.7% completely agreeing and 44.4% agreeing. Furthermore, 85.2% of respondents indicated they would often use AI in the future. Finally, 75.3% expressed willingness to recommend AI to their friends, highlighting its perceived value beyond personal use.

Additionally, the results illustrates students' perceptions of how AI can assist in programming learning, allowing for multiple selections. Providing examples of programming code was the most frequently chosen option (27%), this was closely followed by answering programming questions (24%). Explaining programming concepts (19%) and correcting programming code (17%) were also commonly reported, while 14% selected offering learning advice and resources.

5.2.1 Benefits of Using AI in Programming Learning. Our analysis of students' open-ended responses revealed several perceived benefits when using AI tools such as ChatGPT for programming learning.

Personalized and Accessible Support: Students emphasized the accessibility and personalized nature of AI support. They appreciated being able to ask questions "anytime and anywhere," often in conversational or vague terms that

would be difficult to phrase in a search engine. One respondent noted that “even if I don’t know the right technical words, I can still get a useful answer,” while another highlighted the feeling of “one-to-one support tailored to my needs.”

Error Correction: One of the most frequently mentioned benefits is the ability of AI to quickly identify and correct errors in students’ code. Rather than spending excessive time debugging on their own, students can receive instant feedback, which both saves time and reduces frustration, allowed them to resolve doubts without waiting for help from instructors or peers. As one respondent described, “when I have a question, it solves it immediately,” while another emphasized that they could “move smoothly to the next task without being stuck.” Such immediacy was often framed as a key factor in keeping up learning momentum.

Concept and Code Explanation: Another recurring advantage is that AI provide clear explanations of programming concepts and code logic. Many respondents reported that it could explain programming concepts in simpler terms, often accompanied by concrete usage examples. For instance, students described how it “breaks down complex concepts step by step” and “explains why certain functions are used”. Others appreciated how it supplemented their understanding of Python rules they had previously overlooked, “things I couldn’t understand on my own became obvious after the explanation.” Students described situations where they previously struggled to understand why certain code behaved in a particular way, but through explanations, it became much clearer.

Expanding Ideas and Examples: In addition, AI support creative thinking by suggesting alternative solutions, offering additional examples, and even suggesting new exercises. This encourages students to explore different approaches rather than relying on a single solution. Several participants reported that AI “not only answered but also proposed alternative methods”, while another reflected, “I thought it would just show me one way to do things, but it pointed out where I was wrong and suggested other methods, which made it easier to understand”. Others noted that it provided not only code but also practical examples of function usage, helping them apply knowledge more flexibly.

5.2.2 Concerns of Using AI in Programming Learning.

Accuracy and Reliability Issues: Students repeatedly pointed out that AI responses are not guaranteed to be correct: “You don’t know if that is the correct answer or not because ChatGPT sometimes gives false answers” and “It may not always be accurate and can create mistakes in its responses.” Some noted that the model sometimes “answers based on the question I asked before which has nothing to do with my new questions,” causing confusion. Others mentioned that errors can be subtle and difficult to detect: “If we just prompt the question without thinking first, we wouldn’t learn anything and might absorb wrong information without noticing.”

Overreliance and Reduced Thinking: Another concern is over-dependence, many worried that “getting the answer too quickly makes me stop thinking for myself” or “I just copy the code without reading it, so I don’t gain any knowledge.” Others echoed this concern with remarks like “ChatGPT is too useful, it makes me lazy,” and “If I rely on ChatGPT all the time, my problem-solving ability will decline.” Several stressed that the danger lies not in occasional use but in the habit of outsourcing every step: “I rely too much and my thinking power decreases” and “If you always ask ChatGPT, you stop trying to solve the problem on your own.”

Lack of Context: Students also highlighted the model’s inability to account for course-specific constraints. For example, one noted, “it cannot know what the topic we are learning in class, so it might give answers that students

have not yet encountered.” Others commented that the answers sometimes diverged from what the instructor expected: “Sometimes the direction from AI is different from the teacher’s.”

5.2.3 *Suggestions for Improving AI’ Support in Programming Learning.*

Accuracy and Reliability of Responses: The most consistent demand from students was for more dependable and auditable responses. Participants repeatedly noted that while AI can produce fluent answers, its tendency to “sound confident while being wrong” undermines trust. Typical suggestions included “quantify answer accuracy” and “flag the parts that might be wrong”. Others wanted it to fact-check before responding—“It should check facts before it provides information”, and to continuously update knowledge bases to avoid outdated examples.

In programming contexts, students explicitly asked for verification by execution: “Automatically run generated code in a sandbox and use the results for feedback learning.” Transparency was also emphasized: “Show the reasoning process that led to the answer.” In essence, students are asking not only for accurate output, but for mechanisms of accountability, so that they can gauge whether an answer can be safely trusted.

Instructional Scaffolding and Learning Process Support: Beyond correctness, students wanted AI to act as a learning companion rather than a shortcut to solutions. Many stressed that the system should resist giving full answers too quickly, instead adopting a scaffolded approach: “Instead of the complete answer, provide hints step by step.” Several requested adaptive questioning to check comprehension—“Ask just-right questions to test if I actually understand.”

This also extended to encouraging productive struggle. For instance: “Don’t show the answer right away; push me to think first.” Some even suggested structural restrictions: “Make students go through the reasoning process before they can see the answer” or “Limit how many times we can ask questions so that we think more on our own.”

Students also wanted multiple solution paths and richer learning aids. Comments included: “Don’t just give one example—show different patterns too” and “Explain with detailed charts or diagrams.” These requests suggest that students are envisioning AI as interactive tutor—not merely delivering answers, but scaffolding reasoning, broadening exposure to alternative approaches to deepen understanding.

Context Awareness, Personalization, and User Experience: Another cluster of requests focused on context sensitivity and seamless integration into learning environments. Students expressed frustration when AI produced generic or irrelevant answers and asked for alignment with class materials: “I hope it provides code using the methods the teacher taught us.” Others emphasized focus and stability: “It should stay on topic to the question being asked.”

There was also a call for better prompting support and proactive guidance, e.g., “Help me phrase my question so I can get the answer I want” or “Suggest options without me having to ask first.” In terms of tools, integration with existing workflows was highlighted: “Embed it into IDE so errors and fixes appear directly” and “Allow access to my local files.”

Finally, students imagined multimodal and accessible interfaces: “Explain not only with words but also with illustrations,” and “Add voice commands or visual aids.” Together, these comments point toward a vision of AI as personalized, contextually grounded, and embedded in the actual programming environment, rather than a detached Q&A tool.

6 Educator Survey (RQ3)

We surveyed eleven computing educators from Japan and China to understand their perspectives on using GenAI in programming education. The survey began by exploring the educators’ backgrounds and their current challenges and strategies, especially around students’ utilization of AI-based coding tools. The survey items are provided in Appendix C.

All participants were actively involved in teaching computing courses as well as in computer science education research. Their teaching experience varied: two had more than ten years, two had 3–5 years, three had 1–3 years, and four had less than one year. With respect to teaching level, seven participants taught only undergraduate courses, two taught only graduate courses, and two were engaged in both. Regarding course types, three taught only introductory programming, four taught only data science or machine learning, three taught both introductory programming and data science/ML, and one reported teaching all three types (introductory programming, data science/ML, and PBL). As for programming languages, most participants reported teaching with Python, while three also used C/C++, and one reported Java.

6.1 AI-Policy and General Perceptions

Educators we surveyed generally reported either “mostly open” or “allowed with conditions” policies toward students’ use of AI. Under conditional policies, students were encouraged to consult AI but were explicitly prohibited from directly generating or submitting complete solutions or code. Some instructors required evidence of prior problem-solving effort before any code generation—“If they want to generate code, they need to show that they have done sufficient work towards solving the problem, e.g., decomposing a problem to sub-problems and using AI to generate code for sub-problems”—and others stressed explicit disclosure/acknowledgment of AI use. By contrast, more open policies emphasized AI’s value for personalized learning and for fostering student responsibility and autonomy, particularly at the graduate level. Across both groups, common restrictions included prohibiting AI use during examinations and discouraging reliance on AI for final answers. Overall, educators expressed generally positive views on the role of AI assistants in student learning: 80% of the ten respondents either strongly agreed or agreed that AI assistants help students learn programming, while the remaining 20% were neutral.

6.2 Student Input and Context

6.2.1 Perceived Learning Value of Questions. To further illustrate educators’ perceptions of different types of student prompts, Figure 7(a) summarizes their ratings of the learning value associated with each category across the *planning*, *monitoring*, and *evaluation* phases. Educators consistently associated higher learning value with questions that promote sense-making—such as clarifying concepts, diagnosing errors, and articulating reasoning—than with those that primarily generate complete solutions. Within *planning*, conceptual questions were rated as most beneficial, followed by example requests and implementation questions, and code generation was rated lowest. During *monitoring*, educators valued error interpretation and code correction, whereas code verification elicited more divided views. In *evaluation*, code explanation was recognized as pedagogically useful, and code optimization was often viewed with skepticism, likely reflecting its limited applicability within typical coursework. Collectively, educators believe that AI support and instructional practice should emphasize prompts that foster conceptual framing, diagnostic reasoning, and iterative refinement, rather than encouraging direct solution generation.

6.2.2 Input Formats for AI Use. When asked which input format better promotes high-quality thinking when students use AI, most educators favored structured forms designed by instructors or experts, while a smaller number supported free-form inputs. Proponents of structured formats emphasized that such scaffolds can help students articulate their understanding, avoid incomplete or misleading prompts, and improve programming skills. One educator noted that “many students struggle to provide clear, fully contextualized prompts, which often yield misleading outputs. Well-designed, structured input forms (e.g., templates or checklists) can scaffold better prompts and improve comprehension.”

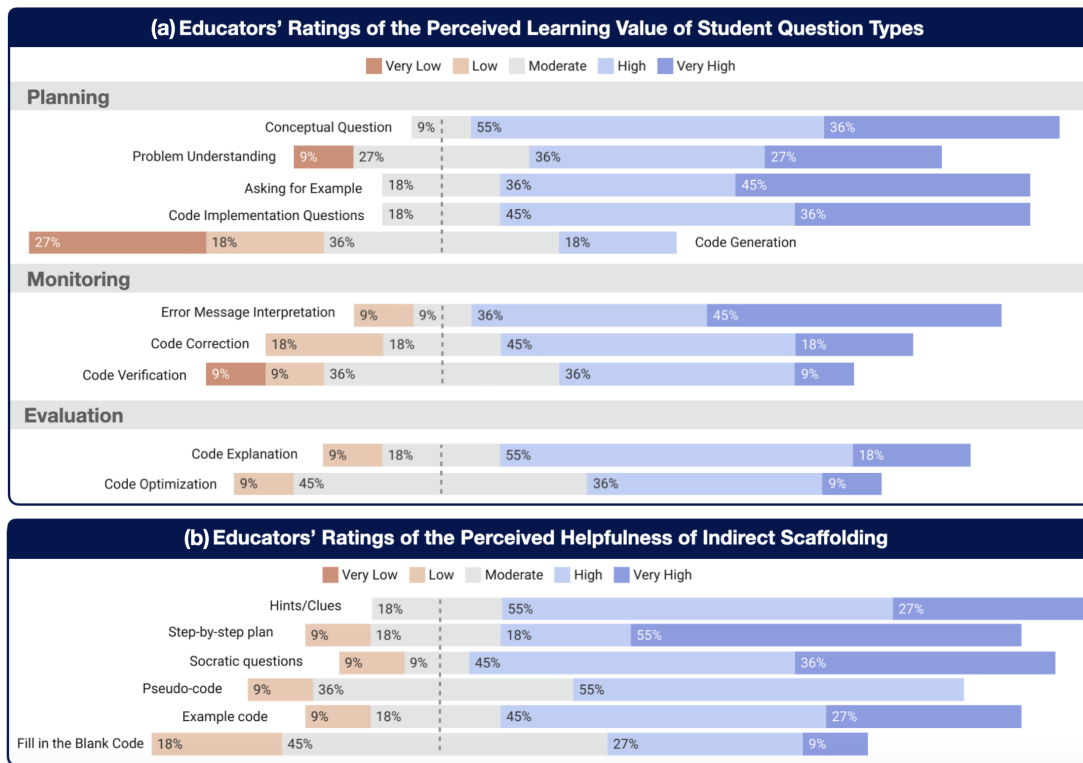


Fig. 7. Educators' ratings of (a) different student question types, and (b) the Perceived Helpfulness of Indirect Scaffolding regarding their perceived learning value.

Others highlighted that structured formats are particularly important for students with weak programming backgrounds, as they might otherwise resort to directly asking AI for answers without sufficient reasoning. In contrast, educators who preferred free-form inputs pointed to accessibility and student autonomy. They argued that free-form interactions encourage more frequent questioning and that “asking AI effectively is also an important skill.” One educator noted that because free-form input mirrors real-world tools such as ChatGPT, it can lower barriers and promote engagement, even if the questions vary in quality.

6.3 AI Responses and Scaffolding

6.3.1 Educator Preferences for Response Scaffolding. All educators consistently agreed that AI's responses should be delivered in the form of indirect scaffolding rather than providing direct runnable code. This shared preference reflects a common belief that scaffolding supports deeper engagement and prevents students from bypassing critical reasoning processes.

Building on this overall preference for indirect scaffolding, Figure 7(b) presents educators' ratings of different forms of such scaffolding. Hints/clues, along with step-by-step plans and Socratic questioning, were rated most helpful, highlighting the value of guidance that supports reasoning and incremental problem solving. In contrast, more code-oriented scaffolds such as pseudo-code and example code received mixed evaluations, while fill-in-the-blank code

was rated least helpful, reflecting skepticism about its pedagogical value. Overall, educators favored scaffolds that provide high-level guidance and reasoning support, rather than more detailed code fragments that risk reducing student engagement in problem solving.

6.3.2 Boundaries and Control of Scaffolding Levels. When asked about the maximum level of scaffolding they would allow, educators’ responses varied, though most favored limiting support to example code or fill-in-the-blank code, with fewer supporting step-by-step plans or pseudocode as the upper bound. Their reasoning highlighted the importance of balancing guidance with student autonomy: educators emphasized that “allowing students to write code themselves using examples ensures their mastery of programming” and that overly complete assistance risks “fostering over-reliance.” At the same time, many acknowledged that beginner students need lighter cognitive burdens and more structured supports, making options such as fill-in-the-blank code appropriate for entry-level contexts.

Across these discussions, there was strong agreement that instructors should play a central role in setting boundaries, while students should retain some flexibility within those boundaries. As one educator explained, “instructors may know better what’s good in general, but students want to personalize the setting to some extent to suit their learning styles.” Others stressed the risk of unregulated autonomy, noting that “if students are given full freedom, they will often gravitate toward the maximum help option, undermining deliberate practice.” A minority suggested more dynamic arrangements, such as AI proposing scaffolding levels with student confirmation, but the overall consensus was that instructor-led constraints are essential to align AI use with pedagogical goals.

6.4 Perceived Benefits, Risks, and Design Considerations

Educators identified several important advantages of using AI assistants in programming education. They emphasized that AI provides timely and personalized support, enabling students to ask questions freely without fear of judgment, receive guidance outside class hours, and access diverse solutions beyond what instructors might offer. AI was also valued for lowering barriers to practice, sustaining motivation, and helping tailor learning to students’ varied backgrounds.

At the same time, educators consistently highlighted drawbacks and risks, particularly the danger of students bypassing critical reasoning by directly generating answers. Concerns included over-reliance, diminished opportunities to make and learn from mistakes, weakened problem-solving and critical thinking skills, and the possibility of AI outputs being misaligned with course goals or overly advanced for novices.

To better promote high-quality thinking, educators suggested that AI should provide scaffolding rather than final solutions—such as hints, step-by-step prompts, Socratic questioning, and personalized guidance based on student context. They stressed the importance of balancing support and fading to encourage independence, and several proposed that AI should actively ask students questions to stimulate reflection rather than supplying complete answers.

To better promote high-quality thinking, educators stressed that AI should provide scaffolding rather than complete solutions, with scaffolds gradually fading to encourage independence. Looking ahead, they envisioned student-facing features such as adjustable response levels, context-aware feedback, and long-term progress tracking, alongside instructor-facing tools including dashboards to monitor interactions, summaries of student queries, and customizable controls over the scope of AI assistance.

7 Design Implications for AI-Powered Coding Assistants with Metacognitive Support (RQ4)

Drawing on our analysis and survey data from students and educators, we derive design implications for AI-powered coding assistants that improve users’ metacognition through support strategies—for example, scaffolds for planning

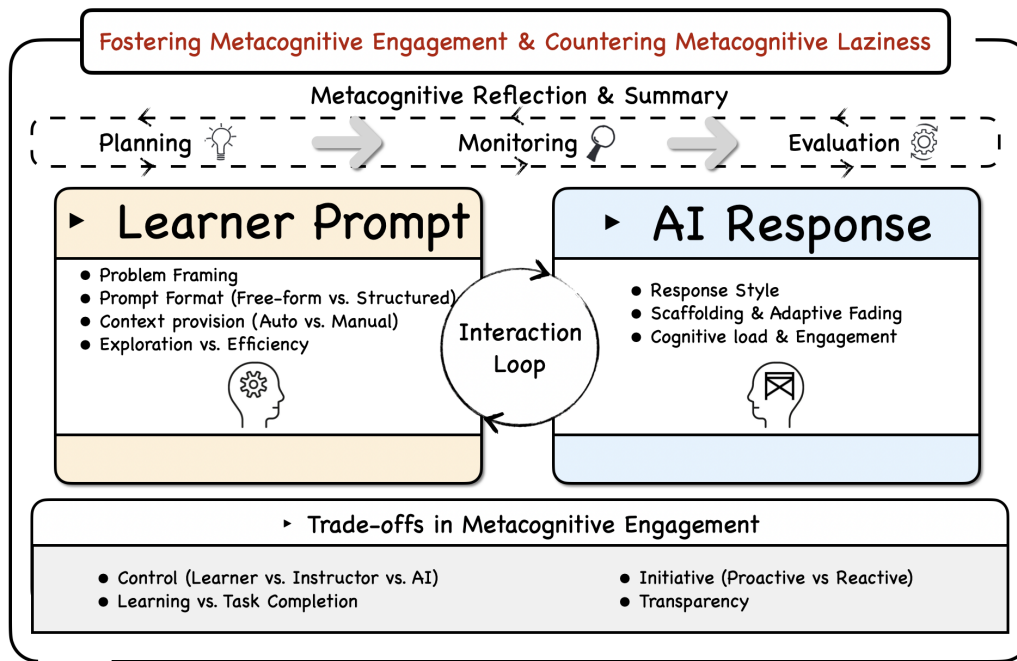


Fig. 8. A Conceptual Design Framework for Metacognition-Supportive AI Programming Assistants.

and evaluation—that can be embedded in GenAI systems. We emphasize that different phases of metacognition entail distinct challenges, risks, and trade-offs, and we highlight opportunities to embed scaffolding that promotes deeper learning rather than superficial task completion. Specifically, as shown in Figure 8, we structure these implications around student–AI help-seeking interactions: (1) the metacognitive phase of AI engagement, (2) the formulation of a prompt, (3) types and scaffolding of AI responses, and (4) navigating trade-offs in metacognitive engagement.

7.1 The Metacognitive Phase of AI Engagement

Our analysis revealed that most student–AI exchanges were confined to a single metacognitive phase, predominantly monitoring, rather than traversing a complete planning-monitoring-evaluation cycle. This fragmentation stems from multiple factors: while some metacognitive activities may occur outside the chat interface (e.g., peer discussions, reviewing notes, consulting learning materials), learners may be overly reliant on AI for immediate assistance rather than engaging in more structured and systematic learning practices [60]. Moreover, they often lack awareness of the full metacognitive cycle or inadvertently outsource cognitive work to the AI—"metacognitive laziness".

To address this, we propose that metacognition-supportive AI assistants should both cultivate learners' metacognitive awareness and guide them through complete metacognitive cycles. While recent research has explored explicit phase scaffolding approaches—such as asking learners to select a self-regulated learning phase before querying or providing phase-specific hint categories [31, 43]—future systems should move toward more integrated conversational designs. These designs would maintain coherence across turns, enabling planning-oriented exchanges to naturally transition into monitoring and ultimately culminate in evaluation phases. In the planning phase, AI can prompt students to explicitly define or restate the problem in their own words, clarify task goals, and specify input–output expectations,

thereby fostering early articulation of intent and externalizing tacit assumptions. During monitoring, rather than directly supplying corrections, the system can encourage learners to predict or explain potential errors, compare alternative solutions, or reflect on discrepancies before receiving hints or corrections, ensuring that error interpretation precedes resolution. The evaluation phase should extend beyond correctness checking to broader transfer of learning, such as asking whether a solution generalizes to new tasks, prompting self-explanation or paraphrasing of the final approach, or encouraging optimization by checking boundary cases, efficiency, and readability. Finally, AI assistants could generate reflective summaries of the interaction that highlight key strategies, phases traversed, and moments of difficulty, prompting learners to step back from the task and recognize patterns in their thinking. Strategic use of brief self-explanation prompts can serve as bridges between phases, counteracting metacognitive “laziness” while keeping learners actively engaged throughout the cycle. Together, these design elements situate AI not merely as a problem-solver, but as a partner in metacognitive regulation.

7.2 The Formulation of a Prompt

Formulating effective prompts is one of the most critical—and metacognitively demanding—activities in AI-assisted programming. While free-form input, common in chat-based models, feels approachable, our analysis shows that many novices struggle to craft clear prompts. Their inputs are often vague and lack context, leading to unhelpful responses. This reflects a genuine challenge: for beginners, articulating precise questions is not trivial [36].

Some researchers have argued for systems that automatically gather context—for instance, IDE-style plugins or course-specific AI assistants using retrieval-augmented generation (RAG)—so that models can deliver more well-aligned responses [28]. Yet unlike productivity-oriented programming assistants, educational tools must prioritize supporting learning over simply producing efficient solutions. In this view, structured prompting plays a crucial role. Many educators emphasized that structured prompts foster active engagement, critical thinking, and intentional inquiry. Without scaffolding, however, novices often struggle to frame effective questions and instead fall back on requesting direct answers from AI.

In addition, systems can scaffold prompt crafting through strategies such as task decomposition and prompt chaining—breaking a complex problem into smaller, sequential sub-tasks where outputs from one step inform the next [57]. These strategies not only improve model performance but also help learners clarify goals and refine their reasoning. Complementary feedforward design can further reduce unproductive trial-and-error by signaling when a vague prompt is unlikely to succeed. After the user specifies a task, the system can generate a step-by-step plan for completion, supported by proactive prompting and concrete examples, while still allowing students to skip decomposition to avoid unnecessary cognitive load. For example, when a student requests “write a Python function to sort numbers,” the system might guide them through clarifying prompts such as: What will the input look like (a list of integers, floats, or strings)? What should the output be (a new sorted list, or the same list modified in place)? Should the sorting be ascending or descending? By helping students articulate precise requirements before code is generated, the system transforms prompt formulation into an intentional learning process rather than a shortcut to solutions.

7.3 Types and Scaffolding of AI Responses

After formulating a question, the next design consideration is the degree of control over the type and directness of assistance. Both students and educators in our study emphasized the importance of indirect scaffolding, which promotes active engagement rather than passive answer-seeking.

Among common forms of assistance such as hints, step-by-step plans, socratic questioning, pseudo-code, example code, and fill-in-the-blank templates, educators generally preferred hints and step-by-step guidance, as it both reflects learners' current progress and reinforces best practices. Socratic questioning, as suggested in prior work [24], was seen as a way to foster independent problem-solving and encourage critical thinking by transforming code generation requests into a series of reflective prompts. Pseudo-code is already used in some systems, yet several teachers worried it may confuse novices who struggle to bridge abstract logic with executable syntax. Example code remains a crucial learning resource [28, 33], but educational AI tools should carefully differentiate between illustrative examples and direct solutions. Fill-in-the-blank code represents another intermediate scaffold, encouraging students to complete missing logic themselves. However, both forms of support are still rarely implemented in current AI assistant tools, and their pedagogical potential warrants further investigation. Importantly, different types of tasks and learner profiles may call for different scaffolding forms. Learners may also exhibit varying preferences depending on factors such as task difficulty, task type, and their stage of learning. Understanding how these dimensions interact remains an open area for future empirical research.

Another dimension concerns the level of help. The key trade-off is providing sufficient scaffolding to promote critical learning while avoiding over-guidance that diminishes autonomy [28]. If support is too direct and learners have complete freedom to choose minimal scaffolding, they may bypass essential learning opportunities within their zone of proximal development, undermining skill growth and self-efficacy. Conversely, overly indirect feedback risks leaving students frustrated and unsupported. To balance these tensions, future systems might adopt adaptive fading strategies.

7.4 Navigating Trade-offs in Metacognitive Engagement

Designing AI systems to support metacognitive engagement requires navigating several trade-offs. While such tools can scaffold planning, monitoring, and evaluation, their effectiveness depends on how agency, visibility, and control are distributed across learners, instructors, and the AI itself. Below we discuss three key considerations: who initiates questions, how transparent the process should be, and who ultimately controls the level of assistance.

One important trade-off concerns the locus of initiative in interactions: should the assistant proactively generate questions and prompts to scaffold metacognition, or should it remain reactive and wait for learners to initiate? Proactive designs may reduce metacognitive inertia and help novices who struggle to ask effective questions, while reactive designs preserve learner autonomy and prevent over-direction. A balanced approach may involve adaptive prompting, where the system intervenes only when learners show signs of confusion, stagnation, or superficial engagement.

Transparency is another key factor for sustaining metacognitive engagement. Learners often benefit from visibility into their own progress and the nature of the scaffolding they receive. For example, an AI assistant could provide a monitoring dashboard that tracks which phases of metacognition (planning, monitoring, evaluation) have been engaged, highlights which types of support (hints, examples, explanations) were requested, and visualizes shifts in the level of assistance. Such transparency can help students reflect on their learning patterns, give instructors oversight to adjust pedagogical strategies, and increase trust in the AI system.

Beyond initiative and transparency, a third dimension concerns who controls the level of assistance. Should the AI system make adaptive decisions on its own, should learners be able to fully direct their experience, or should instructors retain authority over the process? Recent work has shown that learners with higher self-efficacy often prefer greater autonomy in choosing their level of support, suggesting that AI learning tools should be context-aware and dynamically adjust the range of control available to learners. In contrast, preferences for the type of help (e.g., hints vs. examples) may be more closely tied to individual style and taste [24]. A promising direction is to design customizable options

that allow learners to calibrate the AI’s support to their needs while still aligning with instructional goals. Future AI assistants could adopt shared and adaptive control rather than relying on fixed settings. Instructors might define the permissible range of support while learners exercise autonomy within those bounds. Such adaptive approaches not only balance autonomy and guidance but also highlight an important design principle: control of scaffolding should be understood as a dynamic, shared process—shaped by learners, instructors, and context—rather than a static parameter.

Taken together, these four dimensions outline a design space for metacognition-supportive AI in programming education. Supporting learners requires attention not only to the phases of metacognitive engagement (planning, monitoring, evaluation), but also to the formulation of prompts that shape how learners articulate their intent, the types and scaffolding of responses that balance directness with inquiry, and the trade-offs in agency, transparency, and control that govern how assistance unfolds in practice. By connecting these dimensions, we emphasize that AI should not be conceived merely as a source of answers, but as a partner that helps learners navigate metacognitive cycles, develop prompting strategies, and engage with adaptive scaffolding. This integrated perspective highlights opportunities for future research to investigate how these layers interact, and how AI can be designed to sustain metacognitive growth across diverse learners and contexts.

8 Limitations and Future Work

This study has several limitations. First, our findings derive from a single Python programming course at one institution, which constrains the breadth of applicability. Validation across courses in other languages (such as C or Java), as well as across different institutional contexts and cultural environments, will be important for assessing the robustness of these results. Also, the nature of students’ interactions was closely tied to the capabilities of the LLM models available in each year. The GPT-3.5-turbo model used in 2023 may have constrained the quality of support, shaping student perceptions differently from later cohorts who engaged with GPT-4 and GPT-4o. Beyond model versioning, variation in performance also emerges across different programming tasks. Moreover, this study examined only student–LLM interactions, however students may express their critical thinking in ways that are not captured in chatbot logs. In actual classroom contexts, learning unfolds across multiple modalities, such as seeking help from peers or instructors, working within IDEs, consulting learning materials, or referring back to prior exercise solutions. These additional dimensions may play a crucial role in shaping students’ metacognitive processes. Future research should therefore move beyond a single-modal focus by incorporating richer forms of evidence, including classroom video recordings, eye-tracking, and screen-capture data, as well as think-aloud protocols. Such approaches would provide a more comprehensive account of the strategies students adopt and the processes through which they regulate their learning when engaging with AI-powered tools.

9 Conclusion

This paper investigated how students engage with AI in programming education, with a particular focus on metacognition processes. By analyzing more than 10,000 dialogue interactions, we identified how students formulated prompts, received feedback, and navigated planning, monitoring, and evaluation processes with AI support. In addition, we complemented this analysis with student surveys that captured learners’ perceptions and experiences, as well as educator surveys that offered insights into the future role of AI-powered tools in programming education. Synthesizing evidence across these sources, we derived high-level design considerations, highlighting key trade-offs that characterize the emerging design space of educational AI tools. We hope that these findings and design implications will help guide

the development of future LLM-powered coding assistants that support not only problem solving, but also deeper metacognitive engagement and self-regulated learning.

References

- [1] Sara Amani, Lance White, Trini Balart, Laksha Arora, Kristi J Shryock, Kelly Brumbelow, and Karan L Watson. 2023. Generative AI perceptions: A survey to measure the perceptions of faculty, staff, and students on generative AI tools in academia. *arXiv preprint arXiv:2304.14415* (2023).
- [2] Christos-Nikolaos Anagnostopoulos. 2023. ChatGPT impacts in programming education: A recent literature overview that debates ChatGTP responses. *arXiv preprint arXiv:2309.12348* (2023).
- [3] 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.
- [4] Andrea J Bingham and Patricia Witkowsky. 2021. Deductive and inductive approaches to qualitative data analysis. *Analyzing and interpreting qualitative data: After the interview* (2021), 133–146.
- [5] Som Biswas. 2023. Role of ChatGPT in Computer Programming.: ChatGPT in Computer Programming. *Mesopotamian Journal of Computer Science* 2023 (2023), 8–16.
- [6] Joel Brandt, Philip J Guo, Joel Lewenstein, Mira Dontcheva, and Scott R Klemmer. 2009. Two studies of opportunistic programming: interleaving web foraging, learning, and writing code. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. 1589–1598.
- [7] Cecilia Ka Yuk Chan and Katherine KW Lee. 2023. The AI generation gap: Are Gen Z students more interested in adopting generative AI such as ChatGPT in teaching and learning than their Gen X and millennial generation teachers? *Smart learning environments* 10, 1 (2023), 60.
- [8] Angxuan Chen, Mengtong Xiang, Junyi Zhou, Jiyou Jia, Junjie Shang, Xinyu Li, Dragan Gašević, and Yizhou Fan. 2025. Unpacking help-seeking process through multimodal learning analytics: A comparative study of ChatGPT vs Human expert. *Computers & Education* 226 (2025), 105198.
- [9] Gary Cheng, Di Zou, Haoran Xie, and Fu Lee Wang. 2024. Exploring differences in self-regulated learning strategy use between high-and low-performing students in introductory programming: An analysis of eye-tracking and retrospective think-aloud data from program comprehension. *Computers & Education* 208 (2024), 104948.
- [10] Heeryung Choi, Jelena Jovanovic, Oleksandra Poquet, Christopher Brooks, Srećko Joksimović, and Joseph Jay Williams. 2023. The benefit of reflection prompts for encouraging learning with hints in an online programming course. *The Internet and Higher Education* 58 (2023), 100903.
- [11] Ali Darvishi, Hassan Khosravi, Shazia Sadiq, Dragan Gašević, and George Siemens. 2024. Impact of AI assistance on student agency. *Computers & Education* 210 (2024), 104967.
- [12] Paul Denny, Viraj Kumar, and Nasser Giacaman. 2023. Conversing with copilot: Exploring prompt engineering for solving cs1 problems using natural language. In *Proceedings of the 54th ACM Technical Symposium on Computer Science Education V. 1*. 1136–1142.
- [13] Paul Denny, Juho Leinonen, James Prather, Andrew Luxton-Reilly, Thezyrie Amarouche, Brett A Becker, and Brent N Reeves. 2023. Promptly: Using Prompt Problems to Teach Learners How to Effectively Utilize AI Code Generators. *arXiv preprint arXiv:2307.16364* (2023).
- [14] 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. 2023. Computing Education in the Era of Generative AI. *arXiv preprint arXiv:2306.02608* (2023).
- [15] A Ebrahimi, D Kopec, and C Schweikert. 2006. Taxonomy of novice programming error patterns with plan, web, and object solutions. *Comput. Surveys* 38, 2 (2006), 1–24.
- [16] Yizhou Fan, Luzhen Tang, Huixiao Le, Kejie Shen, Shufang Tan, Yueying Zhao, Yuan Shen, Xinyu Li, and Dragan Gašević. 2025. Beware of metacognitive laziness: Effects of generative artificial intelligence on learning motivation, processes, and performance. *British Journal of Educational Technology* 56, 2 (2025), 489–530.
- [17] 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* (Virtual Event, Australia) (*ACE '22*). Association for Computing Machinery, New York, NY, USA, 10–19. <https://doi.org/10.1145/3511861.3511863>
- [18] 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.
- [19] Zhikai Gao, Sarah Heckman, and Collin Lynch. 2022. Who Uses Office Hours? A Comparison of In-Person and Virtual Office Hours Utilization. In *Proceedings of the 53rd ACM Technical Symposium on Computer Science Education - Volume 1* (Providence, RI, USA) (*SIGCSE 2022*). Association for Computing Machinery, New York, NY, USA, 300–306. <https://doi.org/10.1145/3478431.3499334>
- [20] Qiang Hao and Ruohan Liu. 2025. Towards Integrating Behavior-Driven Development in Mobile Development: An Experience Report. In *Proceedings of the 56th ACM Technical Symposium on Computer Science Education V. 1*. 450–456.
- [21] Niklas Humble, Jonas Boustedt, Hanna Holmgren, Goran Milutinovic, Stefan Seipel, and Ann-Sofie Östberg. 2023. Cheaters or AI-Enhanced Learners: Consequences of ChatGPT for Programming Education. *Electronic Journal of e-Learning* (2023), 00–00.
- [22] Michelle Ichinco and Caitlin Kelleher. 2015. Exploring novice programmer example use. In *2015 IEEE Symposium on Visual Languages and Human-Centric Computing (VL/HCC)*. IEEE, 63–71.

- [23] Fatma Gizem Karaoglan Yilmaz and Ramazan Yilmaz. 2022. Learning analytics intervention improves students' engagement in online learning. *Technology, Knowledge and Learning* 27, 2 (2022), 449–460.
- [24] Priyanka Kargupta, Ishika Agarwal, Dilek Hakkani Tur, and Jiawei Han. 2024. Instruct, Not Assist: LLM-based Multi-Turn Planning and Hierarchical Questioning for Socratic Code Debugging. In *Findings of the Association for Computational Linguistics: EMNLP 2024*. 9475–9495.
- [25] Enkelejda Kasneci, Kathrin Seßler, Stefan Küchemann, Maria Bannert, Daryna Dementieva, Frank Fischer, Urs Gasser, Georg Groh, Stephan Günemann, 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 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*.
- [28] Majeed Kazemitabaar, Runlong Ye, Xiaoning Wang, Austin Zachary Henley, Paul Denny, Michelle Craig, and Tovi Grossman. 2024. Codeaid: Evaluating a classroom deployment of an llm-based programming assistant that balances student and educator needs. In *Proceedings of the 2024 chi conference on human factors in computing systems*. 1–20.
- [29] Sam Lau and Philip J Guo. 2023. From "Ban It Till We Understand It" to "Resistance is Futile": How University Programming Instructors Plan to Adapt as More Students Use AI Code Generation and Explanation Tools such as ChatGPT and GitHub Copilot. In *Proceedings of the 2023 ACM Conference on International Computing Education Research-Volume 1*.
- [30] Juho Leinonen, Paul Denny, Stephen MacNeil, Sami Sarsa, Seth Bernstein, Joanne Kim, Andrew Tran, and Arto Hellas. 2023. Comparing code explanations created by students and large language models. In *Proceedings of the 2023 Conference on Innovation and Technology in Computer Science Education V. 1*. 124–130.
- [31] Huiyong Li and Boxuan Ma. 2025. CodeRunner Agent: Integrating AI Feedback and Self-Regulated Learning to Support Programming Education. In *International Conference on Computers in Education*.
- [32] Mark Liffiton, Brad Sheese, Jaromir Savelka, and Paul Denny. 2023. CodeHelp: Using Large Language Models with Guardrails for Scalable Support in Programming Classes. [arXiv:2308.06921](https://arxiv.org/abs/2308.06921) [cs.CY]
- [33] Dastyni Loksa and Amy J Ko. 2016. The role of self-regulation in programming problem solving process and success. In *Proceedings of the 2016 ACM conference on international computing education research*. 83–91.
- [34] Dastyni Loksa, Lauren Margulieux, Brett A Becker, Michelle Craig, Paul Denny, Raymond Pettit, and James Prather. 2022. Metacognition and self-regulation in programming education: Theories and exemplars of use. *ACM Transactions on Computing Education (TOCE)* 22, 4 (2022), 1–31.
- [35] Boxuan Ma, Li Chen, and Shin'ichi Konomi. 2024. Enhancing programming education with ChatGPT: a case study on student perceptions and interactions in a Python course. In *International Conference on Artificial Intelligence in Education*. Springer, 113–126.
- [36] Boxuan Ma, Li Chen, and Shin'ichi Konomi. 2024. Exploring Student Perception and Interaction using CHATGPT in Programming Education. In *21st International Conference on Cognition and Exploratory Learning in the Digital Age, CELDA 2024*. IADIS Press, 35–42.
- [37] Stephen MacNeil, Andrew Tran, Arto Hellas, Joanne Kim, Sami Sarsa, Paul Denny, Seth Bernstein, and Juho Leinonen. 2023. Experiences from using code explanations generated by large language models in a web software development e-book. In *Proceedings of the 54th ACM Technical Symposium on Computer Science Education V. 1*. 931–937.
- [38] Matthew B Miles and A Michael Huberman. 1994. *Qualitative data analysis: An expanded sourcebook*. sage.
- [39] Engineering National Academies of Sciences, Medicine, et al. 2018. *Assessing and responding to the growth of computer science undergraduate enrollments*. National Academies Press.
- [40] Kimberly A Neuendorf. 2017. *The content analysis guidebook*. sage.
- [41] Maciej Pankiewicz and Ryan S Baker. 2023. Large Language Models (GPT) for automating feedback on programming assignments. *arXiv preprint arXiv:2307.00150* (2023).
- [42] Eunsung Park and Jongpil Cheon. 2025. Exploring Debugging Challenges and Strategies Using Structural Topic Model: A Comparative Analysis of High and Low-Performing Students. *Journal of Educational Computing Research* 62, 8 (2025), 1884–1906.
- [43] Tung Phung, Heeryung Choi, Mengyan Wu, Adish Singla, and Christopher Brooks. 2025. Plan More, Debug Less: Applying Metacognitive Theory to AI-Assisted Programming Education. In *International Conference on Artificial Intelligence in Education*. Springer, 3–17.
- [44] Tung Phung, Victor-Alexandru Pădurean, José Cambronero, Sumit Gulwani, Tobias Kohn, Rupak Majumdar, Adish Singla, and Gustavo Soares. 2023. Generative AI for Programming Education: Benchmarking ChatGPT, GPT-4, and Human Tutors. *International Journal of Management* 21, 2 (2023), 100790.
- [45] Prajish Prasad and Aamod Sane. 2024. A self-regulated learning framework using generative AI and its application in CS educational intervention design. In *Proceedings of the 55th ACM Technical Symposium on Computer Science Education V. 1*. 1070–1076.
- [46] James Prather, Paul Denny, Juho Leinonen, Brett A Becker, Ibrahim Albluwi, Michelle Craig, Hieke Keuning, Natalie Kiesler, Tobias Kohn, Andrew Luxton-Reilly, et al. 2023. The robots are here: Navigating the generative ai revolution in computing education. In *Proceedings of the 2023 working group reports on innovation and technology in computer science education*. 108–159.

- [47] Jaakko Rajala, Jenni Hukkanen, Maria Hartikainen, and Pia Niemelä. 2023. "Call me Kiran" -ChatGPT as a Tutoring Chatbot in a Computer Science Course. In *Proceedings of the 26th International Academic Mindtrek Conference*. 83–94.
- [48] Liam Saliba, Elisa Shioji, Eduardo Oliveira, Shaanan Cohny, and Jianzhong Qi. 2024. Learning with style: Improving student code-style through better automated feedback. In *Proceedings of the 55th ACM Technical Symposium on Computer Science Education V. 1*. 1175–1181.
- [49] 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.
- [50] Jaromir Savelka, Arav Agarwal, Marshall An, Chris Bogart, and Majd Sakr. 2023. Thrilled by Your Progress! Large Language Models (GPT-4) No Longer Struggle to Pass Assessments in Higher Education Programming Courses. In *Proceedings of the 2023 ACM Conference on International Computing Education Research-Volume 1*.
- [51] Abdulhadi Shoufan. 2023. Exploring Students' Perceptions of CHATGPT: Thematic Analysis and Follow-Up Survey. *IEEE Access* (2023).
- [52] Leonardo Silva, António Mendes, Anabela Gomes, and Gabriel Fortes. 2024. What Learning Strategies are Used by Programming Students? A Qualitative Study Grounded on the Self-regulation of Learning Theory. *ACM Transactions on Computing Education* 24, 1 (2024), 1–26.
- [53] Marita Skjuve, Asbjørn Følstad, and Petter Bae Brandtzaeg. 2023. The user experience of ChatGPT: Findings from a questionnaire study of early users. In *Proceedings of the 5th International Conference on Conversational User Interfaces*. 1–10.
- [54] Aaron J. Smith, Kristy Elizabeth Boyer, Jeffrey Forbes, Sarah Heckman, and Ketan Mayer-Patel. 2017. My Digital Hand: A Tool for Scaling Up One-to-One Peer Teaching in Support of Computer Science Learning. In *Proceedings of the 2017 ACM SIGCSE Technical Symposium on Computer Science Education* (Seattle, Washington, USA) (SIGCSE '17). Association for Computing Machinery, New York, NY, USA, 549–554. <https://doi.org/10.1145/3017680.3017800>
- [55] Margaret Smith, Yujie Chen, Rachel Berndtson, Kristen M Burson, and Whitney Griffin. 2017. "Office Hours Are Kind of Weird": Reclaiming a Resource to Foster Student-Faculty Interaction. *InSight: A Journal of Scholarly Teaching* 12 (2017), 14–29.
- [56] Dan Sun, Azzeddine Boudouaia, Chengcong Zhu, and Yan Li. 2024. Would ChatGPT-facilitated programming mode impact college students' programming behaviors, performances, and perceptions? An empirical study. *International Journal of Educational Technology in Higher Education* 21, 1 (2024), 14.
- [57] Lev Tankelevitch, Viktor Kewenig, Auste Simkute, Ava Elizabeth Scott, Advait Sarkar, Abigail Sellen, and Sean Rintel. 2024. The metacognitive demands and opportunities of generative AI. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. 1–24.
- [58] H Tian, W Lu, TO Li, X Tang, SC Cheung, J Klein, and TF Bissyandé. [n. d.]. Is ChatGPT the ultimate programming assistant—how far is it?(2023). *arXiv preprint arXiv:2304.11938* ([n. d.]).
- [59] Ahmed Tlili, Boulus Shehata, Michael Agyemang Adarkwah, Aras Bozkurt, Daniel T Hickey, Ronghuai Huang, and Brighter Agyemang. 2023. What if the devil is my guardian angel: ChatGPT as a case study of using chatbots in education. *Smart Learning Environments* 10, 1 (2023), 15.
- [60] Olga Viberg, Jacqueline Wong, Yael Feldman-Maggor, Nora Dunder, and Carrie Demmans Epp. 2025. Chatting with Code: Exploring LLMs as Learning Partners in Programming Education. In *International Conference on Artificial Intelligence in Education*. Springer, 453–461.
- [61] Ramazan Yilmaz and Fatma Gizem Karaoglan Yilmaz. 2023. Augmented intelligence in programming learning: Examining student views on the use of ChatGPT for programming learning. *Computers in Human Behavior: Artificial Humans* 1, 2 (2023), 100005.
- [62] Zhengdong Zhang, Zihan Dong, Yang Shi, Thomas Price, Noboru Matsuda, and Dongkuan Xu. 2024. Students' perceptions and preferences of generative artificial intelligence feedback for programming. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 23250–23258.

A Thematic Analysis Codebook

This appendix includes the final codebook used for the thematic analysis of students' usage of AI.

Table A1. Thematic Analysis Codebook: **Prompt** - What do students ask, and how do these requests reflect their metacognitive processes?

Phase and Prompt Type		Code Description (Strategy)
Planning		Students seek support in designing initial solution strategies
P1	Conceptual Question (CQ)	Ask questions about programming-related conceptual knowledge
P2	Problem Understanding (PU)	Ask the purpose of the question
P3	Asking for Example (AE)	Request example code for specific concepts or functionalities
P4	Code Implementation Question (CI)	Ask specific questions related to code implementation
P5	Code Generation (CG)	Request the generation of code snippets for specific requirements
Monitoring		Students seek support in identifying and resolving issues during code execution
M1	Error Message Interpretation (EM)	Request interpretation of the meaning and causes of error messages
M2	Code Correction (CC)	Ask for debugging and correction of errors in the code
M3	Code Verification (CV)	Ask for verification of the correctness of code
Evaluation		Students seek support in reflecting on and improving code quality
E1	Code Explanation (CE)	Ask for explanation of the functionality and purpose of the code
E2	Code Optimization (CO)	Request the improvement of existing code based on specific needs

Table A2. Thematic Analysis Codebook: **Response** - How much is the response directly revealing the solution?

Dimensions and Codes		Code Description
How does the response reveal the solution?		<i>How much is the response directly revealing the solution?</i>
S1	Exact Solution Code	Generated the code solution to a question
S2	Step to Fix Semantic Issue	Generated the steps required to fix semantic/logical problems, which usually need additional lines to achieve new functionality
S3	Step to Fix Syntax Issue	Generated the steps required to fix minor syntax issues
S4	Step to Fix External Issue	Generated the steps to fix an issue that is not within the code, but about the compilation or execution
S5	Example Code	Generated a generic example for a function, or a generic implementation
S6	Conceptual Explanation	Provides explanations and clarifications of programming concepts
S7	Code Explanation	Provides explanations and commentary on specific code snippets
How technically correct?		<i>Despite the question, how correct is the response from AI?</i>
Correct		Everything including the answer and its explanation is correct
Incorrect		Any part of the answer or explanation is incorrect
How helpful?		<i>Is the response helpful to students? Does it guide them the right direction based on the provided query? Does it identify their potential issues? or is it completely misleading?</i>
Helpful		Answer that allows the student to take one step further, even if it is not arriving at the final solution
Not Helpful		Answer that does not allow students to progress any further, or it is unrelated

B Student Perception Questionnaire

This appendix lists the pre- and post-questionnaire items assessing students' perceptions of GenAI for programming learning.

Table B1. Pre-questionnaire items.

Questions	Item
Pre1	What is your current level of proficiency in Python programming? • Beginner • Intermediate • Advanced
Pre2	Do you currently use GenAI (e.g., ChatGPT)? • I use them regularly. • I used to use them but not anymore. • I haven't used them but I plan to. • I don't use them and have no plans to.
Pre3	How familiar are you with GenAI (e.g., ChatGPT)? • I can explain them in detail. • I can explain them to some extent. • I have a basic understanding. • I have little understanding. • I have no understanding at all.
Pre4	To what extent do you think using GenAI (e.g., ChatGPT) to learn Python programming is beneficial or detrimental? • I think it is positive. • I rather think it is positive. • I am neutral. • I rather think it is negative. • I think it is negative.
Pre5	In your own words, how would you describe GenAI such as ChatGPT? (Open-ended)

Table B2. Post-questionnaire items.

Questions	Item
Post1	To what extent do you think GenAI (e.g., ChatGPT) are helpful for learning programming? (5-point Likert)
Post2	To what extent are you likely to continue using GenAI (e.g., ChatGPT) for learning programming? (5-point Likert)
Post3	To what extent are you likely to use GenAI (e.g., ChatGPT) frequently in your studies? (5-point Likert)
Post4	To what extent are you likely to recommend GenAI (e.g., ChatGPT) to your peers? (5-point Likert)
Post5	How do you think GenAI (e.g., ChatGPT) can support programming learning? (Multiple answers) • Correct programming code. • Answer programming questions. • Provide examples of programming code. • Offer learning advice and resources. • Explain programming concepts.
Post6	In your own words, how would you describe GenAI such as ChatGPT after actually using them? (Open-ended)
Post7	What do you see as the advantages of using GenAI (e.g., ChatGPT) for learning programming? (Open-ended)
Post8	What do you see as the disadvantages of using GenAI (e.g., ChatGPT) for learning programming? (Open-ended)
Post9	In your view, how could GenAI (e.g., ChatGPT) be improved to better support programming learning? (Open-ended)

C Educator Survey

This appendix lists the educator-survey items.

Table C. Educator survey items.

Questions	Item
A. Background	
A1	Teaching experience • 0-1 years • 1-3 years • 3-5 years • 5-10 years • 10-20 years • Over 20 years
A2	Teaching level • Undergraduate • Graduate
A3	Course type • Intro to Programming • Data Science / Machine Learning
A4	Programming languages used in course • Python • C/C++ • Java
A5	Current AI policy in your course • Prohibited • Allowed with Conditions • Mostly Open
A6	Briefly describe your current AI policy and reason in your course (Open-ended)
A7	Overall, AI assistants help students learn programming in my course (5-point Likert)
B. Input & Context	
B1	Please rate the following student question types by perceived learning value • Conceptual Question (5-point Likert) • Asking for Advice / Learning Resources (5-point Likert) • Problem Understanding (5-point Likert) • Asking for Example (5-point Likert) • Code Implementation Questions (5-point Likert) • Code Generation (5-point Likert) • Error Message Interpretation (5-point Likert) • Code Correction (5-point Likert) • Code Verification (5-point Likert) • Code Explanation (5-point Likert) • Code Optimization (5-point Likert)
B2	When students ask for help, which input format better promotes high-quality thinking? • Structured Form (specific structure or format designed by educators/experts) • Free-form (students can ask questions in a free form manner)
B3	Please briefly explain why you chose this option. (Open-ended)
B4	Should the system automatically gather course-specific context, or require manual input from students? • Automatically gather context • Manual input from students
B5	Please briefly explain why you chose this option. (Open-ended)

Table C. Educator survey items (continued).

Questions	Item
C. Scaffolding & Delivery	
C1	How the system's response should be delivered? <ul style="list-style-type: none"> • Direct solutions (e.g., runnable code or full fix) • Indirect scaffolding (e.g., hints/steps/pseudocode, error interpretation)
C2	Please rate the following indirect scaffolding by perceived helpfulness <ul style="list-style-type: none"> • Hints / Clues (5-point Likert) • Step-by-step Plan (5-point Likert) • Socratic Questions (5-point Likert) • Pseudo-code (5-point Likert) • Example Code (5-point Likert) • Fill-in-the-Blank Code (5-point Likert)
C3	Select the maximum indirect scaffolding style you would allow <ul style="list-style-type: none"> • No AI assistance • Only Hints/Clues • Up to Step-by-step plan • Up to Pseudo-code • Up to Example code • Up to Fill in the Blank Code • Up to Direct runnable code
C4	Please briefly explain why you chose this option (Open-ended)
C5	Who should control the answer's scaffolding style? <ul style="list-style-type: none"> • Instructor controls all • Instructor sets a maximum, students choose within range • Student chooses freely • System suggests and student confirms • System controls all
C6	Please briefly explain why you chose this option (Open-ended)
D. Pros, Cons, and Feature Wishes	
D1	What are the most important advantages of students using AI assistants for programming learning in your course? (Open-ended)
D2	What are the most important drawbacks or risks of students using AI assistants in your course? (Open-ended)
D3	What kind of features of the system do you want to better promote students' high-quality thinking? (Open-ended)
D4	From your teaching perspective, what features would you expect the system to provide? (Open-ended)