

---

# Synthesizing High-Quality Programming Tasks with LLM-based Expert and Student Agents

---

**Manh Hung Nguyen**  
MPI-SWS, Germany  
manguyen@mpi-sws.org

**Victor-Alexandru Pădurean**  
MPI-SWS, Germany  
vpadurea@mpi-sws.org

**Alkis Gotovos**  
MPI-SWS, Germany  
agkotovo@mpi-sws.org

**Sebastian Tschiatschek**  
University of Vienna, Austria  
sebastian.tschiatschek@univie.ac.at

**Adish Singla**  
MPI-SWS, Germany  
adishs@mpi-sws.org

## Abstract

Generative AI is transforming computing education by enabling the automatic generation of personalized content and feedback. We investigate its capabilities in providing high-quality programming tasks to students. Despite promising advancements in task generation, a quality gap remains between AI-generated and expert-created tasks. The AI-generated tasks may not align with target programming concepts, could be incomprehensible to students, or may contain critical issues such as incorrect tests. Existing works often require interventions from human teachers for validation. We address these challenges by introducing PYTASKSYN, a novel synthesis technique that first generates a programming task and then decides whether it meets certain quality criteria to be given to students. The key idea is to break this process into multiple stages performed by expert and student agents simulated using both strong and weaker generative models. Through extensive evaluation, we show that PYTASKSYN significantly improves task quality compared to baseline techniques and showcases the importance of each specialized agent type in our validation pipeline. Additionally, we conducted user studies using our publicly available web application and show that PYTASKSYN can deliver high-quality programming tasks comparable to expert-designed ones while reducing workload and costs, and being more engaging than programming tasks that are available in online resources.

## 1 Introduction

Generative AI is transforming learning and teaching in computing education [1, 2]. Advanced generative models such as OpenAI’s GPT-4o [3] and GitHub Copilot [4] are quickly reshaping both student and teacher experiences. For students, these models can provide personalized educational content [5], provide detailed feedback on their work [6, 7], and serve as pair programmers [8]. For teachers, these models can assist in analyzing and grading student answers [9] and curriculum development [10]. A particularly promising application is their ability to generate tailored educational materials, especially in creating diverse programming exercises that target specific concepts.

Recent works have investigated the use of generative models for generating novel and engaging programming exercises related to a specific theme and targeting specific programming concepts [11–15]. While these initial efforts show promise, AI-generated tasks still fall short of human expert quality due to several issues [9, 12, 14]. For example, the generated programming task can contain incorrect test cases generated as part of the task or it can contain a task description that is not comprehensible [12, 14]. Without an automatic validation mechanism to check these aspects of the

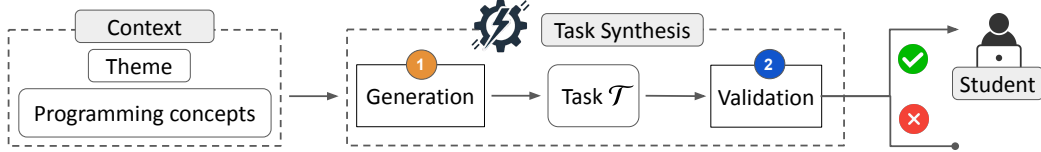


Figure 1: Overview of contextualized programming task synthesis pipeline. Given an input context consisting of a theme and a set of programming concepts, a task  $\mathcal{T}$  is first generated in Stage ① and then validated in Stage ②. Task  $\mathcal{T}$  is assigned to a student only if it meets certain quality criteria.

task, human interventions would be required to validate the task’s quality before they are assigned to students [12, 16]. While one might consider using generative models for task validation, research has shown that they struggle with self-correction [17]. These limitations in single-agent validation motivated our multi-agent based approach.

To address these challenges, we introduce our technique, PYTASKSYN, capable of generating contextualized programming tasks and then deciding whether they meet certain quality criteria. Figure 1 depicts the overview of the pipeline implemented in PYTASKSYN. First, given a context as input, PYTASKSYN asks a simulated agent to generate a programming task composed of a description and a test suite. The generated task then goes through the second stage in the pipeline (i.e., validation), handled by multiple simulated agents with unique roles using strong (GPT-4o [3]) and weaker (GPT-4o-mini [18]) models. This validation stage is designed to provide a quality assurance mechanism and decides whether the generated task can be provided to a student or not. Our multi-agent approach builds upon research that has shown the benefit of collaborative agents [19–21] and used generative models to simulate humans in various roles [7, 22–25]. We demonstrate the efficacy of our technique through extensive evaluation using data collected from prior works [11, 13, 14, 16].

In summary, our main contributions are as follows.

- (1) We highlight the effectiveness of decomposing the programming task synthesis into multiple stages. We introduce PYTASKSYN, a novel technique that leverages generative models to simulate expert, tutor, and student agents, each responsible for specific stages in the synthesis process. (Section 4).<sup>1</sup>
- (2) We conduct extensive evaluation of our technique for Python programming task synthesis, demonstrating significant improvements in the quality of synthesized tasks, while maintaining substantial coverage (Section 5).
- (3) We develop a web application for Python programming task synthesis and conduct two user studies, showing that our synthesized tasks match the quality of expert-created tasks while requiring minimal cost. (Section 6).

## 2 Related Work

**Programming task generation.** With the advancement of generative models, recent studies have demonstrated the potential of using them to create natural language programming tasks [9, 11–14, 16]. Initial research focused on the generation of programming exercises, including sample solutions and test cases [14], while subsequent research efforts have aimed at synthesizing more diverse and contextualized programming exercises [11, 13, 16], as well as exercises suitable for non-English contexts [12]. A parallel line of work focused on generating programming tasks to target specific types of bugs [9], or visual programming tasks [26–28]. However, many of these works highlighted the significant gap between AI-generated tasks and human-crafted tasks [9, 12, 14, 27]. Our work addresses this gap by automating the validation of programming tasks’ quality.

**Improving quality of programming content.** Recent works have focused on enhancing the quality of programming content. For example, generative models were used to generate predicates for testing a generated program [29], or to act as judges for evaluation [30]. As an enhancement, generative models can leverage their own evaluations to iteratively improve the quality of their responses

<sup>1</sup><https://github.com/machine-teaching-group/aied2025-pytasksyn>




<b>Theme:</b> Superheroes		
<b>Programming concepts:</b> Dictionaries , Classes & Objects, Strings, Arithmetic Operators		
<p align="center"><b>Task Description</b></p> <p>Design and implement a Python class ‘Superhero’ that models a simple superhero character using the following guidelines:</p> <ol style="list-style-type: none"> <li>Attributes: <ul style="list-style-type: none"> <li>‘name’ (string) : Name of the superhero</li> <li>‘power’ (string) : A short description of their superpower</li> <li>‘age’ (integer) : Age of the superhero</li> <li>‘world_saving_points’ (integer) : Points representing the superhero’s achievements.</li> </ul> </li> <li>Methods: <ul style="list-style-type: none"> <li>‘__init__(self, name, power, age)’ : This method should initialize a superhero with the provided name, power, and age. The ‘world_saving_points’ should start at 0.</li> <li>‘save_the_day(self, difficulty)’ : This method takes a difficulty level (integer) and increases the ‘world_saving_points’ by two times the difficulty level. If ‘difficulty’ is less than 1, it should not change the points.</li> <li>‘get_description(self)’ : This method returns a string describing the superhero in the format: “{name} possesses the power of {power} and is {age} years old.”</li> </ul> </li> <li>Functions: <ul style="list-style-type: none"> <li>Implement a standalone function ‘top_hero(hero_list)’ that takes a list of ‘Superhero’ objects and returns the name of the superhero with the most ‘world_saving_points’. If there is a tie, return the lexicographically smaller name.</li> </ul> </li> </ol>		
<p align="center"><b>Test suite</b></p> <pre>def test_top_hero():     superheroes = [Superhero("Thor", "thunder god", 1500), Superhero("Hulk", "super                         strength", 35), Superhero("Doctor Strange", "magic", 45)]     superheroes[0].save_the_day(10)     superheroes[1].save_the_day(10)     superheroes[2].save_the_day(12)     assert top_hero(superheroes) == "Doctor Strange"     superheroes[1].save_the_day(4)     assert top_hero(superheroes) == "Doctor Strange" ... (other test cases are omitted for brevity)</pre>		
<p align="center"><b>Validation Result</b></p> <div> <div>  <p>Task description</p> </div> <div>  <p>Test suite</p> </div> <div>  <p>Context relevance</p> </div> </div>		

Figure 2: Example of a contextualized programming task. This task has an **incorrect test** and fails to cover the **“Dictionaries”** concept. Our technique PYTASKSYN validates the quality of this task and abstains from outputting it to students.

[31, 32]. However, prior work has shown that validating the feedback generated by one generative agent using another generative agent has proved to be more effective than using the same agent [7, 24]. Similarly, in our work, we break the task synthesis process into smaller processes and assign them to different agents.

**Generative models as simulated agents.** A recent line of research involves agents assuming different roles and collaborating to achieve a goal. For example, in AutoGen [21], multiple agents can converse, use tools, and incorporate human input. Similar studies show that simulated agents interacting with each other outperform a single agent in reasoning and planning [19, 20]. Several studies have explored LLM-empowered agents for simulating classroom interactions [33], facilitating tutor training [34], and implementing learning-by-teaching environments [35–37]. Our work builds upon these successes by simulating experts, tutors, and students while synthesizing tasks to ensure their quality.

### 3 Problem Setup

We define contextualized programming tasks in Section 3.1, introduce a quality metric in Section 3.2 and technique evaluation metrics in Section 3.3.

### 3.1 Contextualized Programming Tasks

**Context.** We define a context  $\psi$  as a tuple  $\psi = (\psi_{\text{theme}}, \psi_{\text{concepts}})$ , where  $\psi_{\text{theme}}$  represents a theme of interest, and  $\psi_{\text{concepts}}$  denotes a set of target programming concepts for practicing while solving a task. The inclusion of such context aligns with prior works on generating contextualized programming tasks [11, 13, 16].

**Programming task.** We define a programming task  $\mathcal{T}$  as being composed of a task description  $\mathcal{T}_{\text{desc}}$  and a test suite  $\mathcal{T}_{\text{tests}}$ . The task description  $\mathcal{T}_{\text{desc}}$  explains what needs to be accomplished, including the requirements, expected functionality, and constraints. The test suite  $\mathcal{T}_{\text{tests}}$  is a set of test cases used to verify the correctness of a code with respect to  $\mathcal{T}_{\text{desc}}$ . A student is given the task description  $\mathcal{T}_{\text{desc}}$  to solve, while the test suite  $\mathcal{T}_{\text{tests}}$  is kept hidden to verify the student’s code. A student solves task  $\mathcal{T}$  successfully if their code passes all test cases in  $\mathcal{T}_{\text{tests}}$ . Figure 2 shows an example of a contextualized programming task.

### 3.2 Quality of Contextualized Programming Tasks

Ensuring the quality of a contextualized programming task is crucial before giving it to students. It must be relevant to the desired theme and programming concepts while being correct and comprehensible for students to solve. We propose a systematic evaluation process for a human expert to assess the quality of task  $\mathcal{T}$  created for context  $\psi$ . The process begins with the expert formulating a solution for  $\mathcal{T}$ , which is a natural approach to gain a thorough understanding of the task. If the expert cannot formulate a solution, this indicates a fundamental issue with the task, marking it as low-quality. Upon successfully formulating a solution code  $C^*$ , the expert evaluates the correctness of  $\mathcal{T}$  by verifying whether  $\mathcal{T}_{\text{tests}}$  correctly validates  $C^*$ , covering all base and corner cases handled by  $C^*$ . Next, they verify whether task  $\mathcal{T}$  is relevant to  $\psi_{\text{theme}}$  and whether all the programming concepts in  $\psi_{\text{concepts}}$  are required when formulating  $C^*$ . Then, they assess whether the task description  $\mathcal{T}_{\text{desc}}$  provides sufficient information for students to write solutions. We define *Q-Overall* as the overall quality metric, where an expert assigns a final numerical score at the end of the evaluation process.

### 3.3 Technique Evaluation Metrics

To guide our technique development and evaluate its effectiveness, we use two performance metrics: (i) *Coverage*, measuring the percentage of times a technique provides a programming task to the student, and (ii) *Precision*, measuring the percentage of times the programming task provided to the student is of high quality. Our objective is to develop a technique with high precision to ensure well-designed tasks and high coverage rate to provide tasks across diverse contexts.

## 4 Our Technique PYTASKSYN

In this section, we first discuss the motivation and overview of our technique in Section 4.1. Then, we detail each stage in our technique’s pipeline, including task *generation* in Section 4.2 and task *validation* in Section 4.3.

### 4.1 Motivation and Overview

Common validation techniques use a single generative agent to generate both a task and its solution code for consistency checking [14] or to act as a judge for evaluating the output task [30]. However, our initial experiments revealed task quality issues when relying on a single agent (cf. Section 5.4). Previous research show that breaking a complex process into sub-processes and combining the capabilities of multiple agents in a modular manner significantly improves performance [19, 21]. Building on this insight, we implement a multi-agent technique where each agent takes a unique role. Figure 3 presents an overview of our technique consisting of a generation stage and a validation stage. The validation aims to improve precision, but may result in less coverage. To tackle coverage drop, we use an outer loop that repeats synthesis until a task passes validation or a maximum of  $N$  trials is reached; if none pass, no task is given to the student.

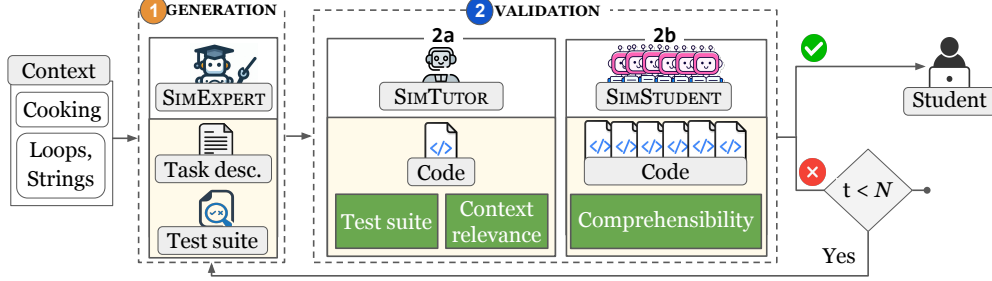


Figure 3: Pipeline of our technique PYTASKSYN. A task candidate is generated by a SIMEXPERT agent given a theme and programming concepts (Stage 1). Then, SIMTUTOR and SIMSTUDENT agents attempt to solve the task before assessing different quality aspects of the generated task candidate (Stage 2 a-b). If the task does not pass the validation stage, our technique retries up to  $N$  times, before ultimately deciding whether or not to assign a task to the student.

## 4.2 Stage 1 - Generation

In this stage, we aim to generate a task  $\mathcal{T}$  composed of a task description  $\mathcal{T}_{\text{desc}}$  and a test suite  $\mathcal{T}_{\text{tests}}$  for a given input context  $\psi$ . For this, our technique uses a simulated expert agent called SIMEXPERT, implemented using the state-of-the-art generative model GPT-4o [3]. Besides asking for the components of a task  $\mathcal{T}$ , we also request a solution code for  $\mathcal{T}$ , drawing inspiration from Chain-of-thought reasoning which has proven to enhance output quality from generative models [38]. We execute this code against  $\mathcal{T}_{\text{tests}}$  as a *generation consistency* check [14]. If it does not pass  $\mathcal{T}_{\text{tests}}$ , task  $\mathcal{T}$  is considered invalid and will not go to the subsequent stages. Figure 4 (Stage 1) shows the structure of the prompt we use for this agent, with a system prompt establishing its role as a programming expert. The user prompt provides the task context ( $\psi_{\text{theme}}, \psi_{\text{concepts}}$ ) and an instruction to generate a programming task. We use the model’s default temperature of 1.0.

## 4.3 Stage 2 - Validation

**Stage 2a - Validation by SIMTUTOR.** In this stage, our goal is to evaluate the test suite and contextual relevance of the task generated in Stage 1. To this end, we implement a SIMTUTOR agent to simulate a human tutor for validating the task. Figure 4 (Stage 2a) shows the prompt we use for this SIMTUTOR agent. It is assigned the role of a tutor through a system prompt, followed by a user prompt containing an input context ( $\psi_{\text{theme}}, \psi_{\text{concepts}}$ ), and a generated task  $\mathcal{T}$ . First, we instruct the SIMTUTOR agent to solve the task by writing a solution code  $\mathcal{C}^*$ . This approach follows the Chain-of-Thought prompting strategy [38], which mimics the evaluating process of a human expert in Section 3.2. Solution  $\mathcal{C}^*$  is then used to validate  $\mathcal{T}_{\text{tests}}$ . Specifically, we verify whether all test cases pass when executing  $\mathcal{C}^*$  and whether every line of  $\mathcal{C}^*$  is covered by  $\mathcal{T}_{\text{tests}}$ . Then, the agent assesses the relevance of the task to the given context. It assigns a score of 1 if task  $\mathcal{T}$  effectively integrates the given theme and programming concepts, and a score of 0 otherwise. We note that if the agent attempts to cheat by writing a code  $\mathcal{C}^*$  that directly uses input-output pairs from  $\mathcal{T}_{\text{tests}}$ , the context relevance score will be 0 as it fails to use the required programming concepts. We use the GPT-4o model with its default temperature of 1.0 for our SIMTUTOR agent.

**Stage 2b - Validation by SIMSTUDENT.** In this stage, we aim to evaluate the comprehensibility of the generated task  $\mathcal{T}$ . Specifically, the task description  $\mathcal{T}_{\text{desc}}$  should provide all the necessary information for a student to write a solution code. To this end, we use a classroom-scale population of SIMSTUDENT agents (20 agents in our experiments) to simulate students’ points of view and obtain their solutions. If the majority of these simulated student agents fail to solve the task,  $\mathcal{T}_{\text{desc}}$  likely lacks clarity or critical information. We consider  $\mathcal{T}_{\text{desc}}$  comprehensible if at least  $\tau$  percent (default at 50%) of the SIMSTUDENT agents successfully solve task  $\mathcal{T}$  given only  $\mathcal{T}_{\text{desc}}$ . To simulate students, we use a system prompt that assigns the agents the role of students, followed by a user prompt containing only the task description (see Figure 4, Stage 2b). We use the GPT-4o-mini [18] model with a default temperature of 1.0, as prior research has shown that weaker models are better suited for simulating students’ perspectives [7].

Stage	Prompt template
1	[System] You are an expert in Python programming. [User] Given a theme of {theme} and a list of programming concepts of {concepts}, generate a Python programming task that requires only the given programming concepts to solve. The task includes a description, a test suite, and a solution program.
2a	[System] You are a tutor in a Python programming course. [User] The following Python programming task was created given a theme of {theme} and a list of programming concepts {concepts}. Task description: {task_description} Test suite: {testsuite} Write a program to solve the task and evaluate the context relevance of the task. The context relevance is 1 if the task is clearly relevant to the given theme and the theme is explicitly used throughout, and all given programming concepts are strictly required to solve the task; 0 otherwise.
2b	[System] You are a student enrolled in a Python programming course. [User] Write a program to solve the task below. Task description: {task_description}

Figure 4: Overview of prompt templates for each stages implemented in PYTASKSYN. {placeholders} are used to include details for concrete scenarios.

## 5 Evaluation

In this section, we present experimental evaluations centering around the following questions: (1) **RQ1**: How does PYTASKSYN perform compared to other existing techniques?; (2) **RQ2**: What are the contributions of different agents in PYTASKSYN?; and (3) **RQ3**: How does PYTASKSYN perform across different contexts? We present our evaluation setup in Section 5.1, 5.2, 5.3, followed by results in Section 5.4. While our evaluations focus on the Python programming language, our technique can be extended to other programming languages.

### 5.1 Context Selection

We evaluate the effectiveness of our technique across varied themes and programming concepts collected from prior works [11, 13, 14, 16]. We select 5 diverse themes and uniformly sample 5 sets of 3 to 5 core Python programming concepts for each theme, resulting in 25 contexts in total (cf. Figure 7 for examples).

### 5.2 Evaluation Procedure

For each sampled context  $\psi$ , we generate  $N = 10$  programming tasks, resulting in a pool of 250 tasks (25 contexts  $\times$  10 tasks). This is done prior to applying different validation techniques in Section 5.3. Two authors, with expertise in Python programming, evaluated the tasks using the Q-Overall metric introduced in Section 3.2 and assigned binary scores (1-High quality/0-Low quality) to each task.<sup>2</sup> To better understand the reasoning behind evaluations, we ask the two annotators to answer three additional Yes(1)/No(0) questions: (i) Is the test suite correct and sufficiently covering relevant cases?; (ii) Does the task accurately reflect the input context?; (iii) Is the task description comprehensible? We obtain a Cohen’s Kappa agreement [39] of 0.8 for Q-Overall metric and at least 0.7 for each additional question, indicating substantial agreement between the annotators. We aggregate the two annotators’ scores to obtain average quality scores for each task. Tasks passed validation of each technique in Section 5.3 are used for computing precision and coverage.

### 5.3 Techniques Evaluated

**Baselines.** First, BASE neither applies a consistency check during the generation stage nor uses any validation mechanisms in the validation stage. Second, GENCONSISTENCY incorporates only

<sup>2</sup>Additionally, our expert evaluation shows that 98% of the 250 generated tasks do not contain programming concepts more advanced than those in the input context.

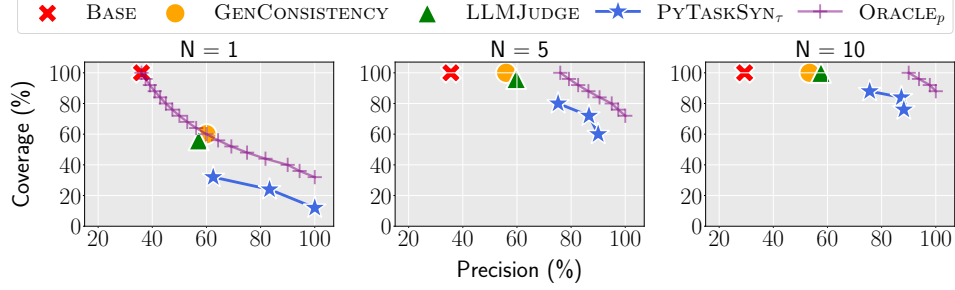


Figure 5: Technique comparison. PYTASKSYN<sub>τ</sub> improves task quality substantially over both baselines GENCONSISTENCY and LLMJUDGE, getting closer to the performance of ORACLE<sub>p</sub>. We vary threshold  $\tau \in \{0\%, 50\%, 100\%\}$  to showcase the tunable precision-coverage trade-off of PYTASKSYN<sub>τ</sub>.

		Q-Overall	Test suite	Context relevance	Comprehensibility
Technique	BASE	0.29	0.42	0.78	0.87
	GENCONSISTENCY	0.53	0.77	0.76	0.85
	LLMJUDGE	0.58	0.80	0.76	0.88
	SIMTUTORSVAL	0.76	0.93	0.91	0.88
	SIMSTUDENTSVAL	0.62	0.83	0.74	0.99
	PYTASKSYN	0.87	0.96	0.92	0.98

Figure 6: Ablation study. We report average Q-Overall scores and answers for three questions answered by experts (cf. Section 5.2). Validation from the simulated tutor agent improves test suite quality and contextual relevance (as shown by SIMTUTORVAL), while validation from simulated student agents enhances task comprehensibility (as shown by SIMSTUDENTVAL). These agents collectively contribute to the overall quality of tasks, as demonstrated by PYTASKSYN.

a consistency check during the generation stage [14]. Third, LLMJUDGE leverages an LLM as a judge [30] to validate a task in the validation stage. We prompt it with input context  $\psi$ , a task  $\mathcal{T}$ , then instruct it to assess the test suite, contextual relevance, and task comprehensibility. It assigns a binary Q-Overall score (1/0) for the task’s overall quality.

**Our Technique and Ablations.** PYTASKSYN involves consistency check in Stage 1 and multiple validation mechanisms in Stage 2. It can be parameterized by  $\tau$ , denoted as PYTASKSYN<sub>τ</sub>, where  $\tau$  represents the threshold for the percentage of SIMSTUDENT agents that solved the task. We implement two ablation variants, each corresponding to a stage of our validation pipeline. SIMTUTORSVAL relies solely on validation from the simulated tutor agent, while SIMSTUDENTVAL uses only validation from simulated student agents. The default value of  $\tau$  is set to 50% for both SIMSTUDENTVAL and PYTASKSYN in our evaluation.

**Oracle.** To evaluate the validation efficacy, we introduce ORACLE<sub>p</sub>, which has access to ground-truth quality of tasks assessed by the human experts. It can select tasks to meet any precision threshold  $p$ , serving as an upper bound.

## 5.4 Results

**RQ1: How does PYTASKSYN perform compared to other techniques?** Figure 5 illustrates the precision and coverage of various techniques. Baseline methods including BASE, GENCONSISTENCY, and LLMJUDGE achieve high coverage but low precision ( $\leq 60\%$ ) due to no or insufficient validation. The subpar precision of LLMJUDGE stems from a single agent’s inability to thoroughly validate all aspects of task quality. Our technique PYTASKSYN<sub>τ</sub> leveraging perspectives from different simulated tutor and student agents improves the quality of synthesized tasks. As  $N$  increases, PYTASKSYN<sub>τ</sub> not only increases coverage but also consistently achieves the highest precision compared to baselines. At  $N = 10$ , PYTASKSYN demonstrates strong performance, achieving a high precision of 87.3% while maintaining substantial coverage at 84.0%. However, there is room for improvement in terms of both precision and coverage when compared to ORACLE<sub>p</sub>. Finally, varying the passing threshold  $\tau$  in PYTASKSYN<sub>τ</sub> reveals a clear precision-coverage trade-off, allowing control based on requirements.



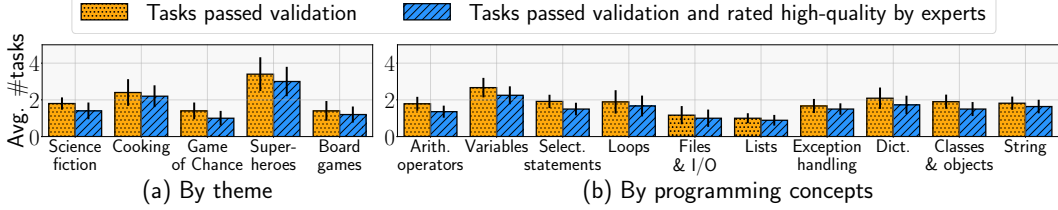


Figure 7: Average number of tasks passed validation of PYTASKSYN out of  $N = 10$  generated tasks, and average number of those rated as high-quality. PYTASKSYN could deliver high-quality tasks for any theme and programming concept.

**RQ2: What are the contributions of different agents in PYTASKSYN?** Figure 6 presents the ablation study analyzing the impact of different agent types on different task quality aspects. Our analysis reveals distinct and complementary contributions from the simulated agents. Simulated tutor agents enhance test suite quality notably and ensure strong context alignment (as shown by SIMTUTORVAL), while simulated student agents substantially boost task comprehensibility (as demonstrated by SIMSTUDENTVAL). Their combination in PYTASKSYN achieves the highest overall quality, with Q-Overall score of 0.87. This demonstrates how our multi-agent validation approach effectively combines different perspectives to ensure comprehensive task quality assessment.

**RQ3: How does PYTASKSYN perform across different contexts?** Figure 7 shows the average number of tasks synthesized by PYTASKSYN as well as how many were rated as high-quality by experts across different themes and programming concepts. We observe that our technique performs well across most themes and programming concepts. Out of  $N = 10$  tasks, 1 to 3 tasks passed validation on average, with the majority rated as high-quality by human experts.

## 6 User Studies Using Our Web Application

We created a public web application (<https://pytasksyn.netlify.app>) where users can request tasks from PYTASKSYN for their chosen themes and programming concepts, write code and debug using our integrated programming environment. We conducted the following two user studies via this web application.

### 6.1 Comparison with Expert-created Tasks and Online Resources

**Setup.** We compared tasks from PYTASKSYN against tasks created by experts and those from online resources. From 5 themes in Section 5.1, we re-sampled 3 to 5 programming concepts to create 5 new contexts. For each, we collected 3 tasks from: (1) an expert, (2) online resources (specifically [geeksforgeeks.org](https://www.geeksforgeeks.org)), and (3) our web application.<sup>3</sup> We recruited 10 volunteer participants, including tutors and graduate students. They are non-native English speakers, with an average age of 28.3. Among them, 5 held Master’s degrees and 5 held Bachelor’s degrees in STEM fields, with an average of 6.8 years of Python experience. We assign each participant 2 random contexts with their corresponding 6 tasks (task sources are undisclosed). After completing each task or reaching the 20-minute limit, they provided feedback using multi-level Likert scales on: theme relevance (1-Yes/0.5-Partial/0-No), programming concepts relevance (1-Yes/0.5-Partial/0-No), task comprehensibility (1-Yes/0-No), difficulty (1-Hard/0.5-Medium/0-Easy), and interestingness (1-Interesting/0.5-Okay/0-Boring).

**Results.** The multi-level feedback provided by participants is mapped to scores in the range  $[0.0, 1.0]$ , with mean scores for each source reported in Figure 8. Additionally, we compute averaged participant success rate, solving time, and time taken to get a task from each source. Our analysis revealed that tasks synthesized by PYTASKSYN achieved quality comparable to expert-designed tasks while requiring significantly less creation time. Moreover, PYTASKSYN has an average task creation cost of just 0.13 USD (APIs cost). Tasks from online resources, while readily available, generally

<sup>3</sup>We made our best effort to find tasks on [geeksforgeeks.org](https://www.geeksforgeeks.org) that covered programming concepts in each context. The presence of specific themes generally cannot be found.



	Theme	Concepts	Compre.	Difficulty	Interest.	Success (%)	Solving time (mins)	Creation time (mins)
Expert	0.95	0.95	0.95	0.42	0.72	84	10.82	25.20
Online resources	0.15	0.82	0.85	0.15	0.50	100	3.26	0
PYTASKSYN	0.98	0.92	0.90	0.20	0.65	100	6.66	1.29

Figure 8: Comparison of tasks from expert, online resources, and PYTASKSYN. We report averaged participant ratings on various task aspects and other statistics. Tasks from PYTASKSYN achieved quality comparable to expert-designed ones while requiring substantially less creation time.

Statistic	Participant ID										Average
	P1	P2	P3	P4	P5	P6	P7	P8	P9	P10	
No. Requests	6	5	6	6	5	5	6	5	5	5	5.40
No. Synthesized Tasks	5	5	5	5	5	5	5	5	5	5	5.00
No. Solved Tasks	5	5	4	4	4	5	5	5	3	5	4.50
Avg. Time to Solve (mins)	4.94	12.74	12.66	13.28	5.46	8.15	10.14	6.92	7.50	6.69	8.85

Figure 9: Statistics of each participant when they were instructed to request and solve tasks synthesized by PYTASKSYN on our web application. Our system successfully synthesized tasks for 92.6% of participant requests.

lacked thematic integration and were rated as less interesting and easy. Expert-created tasks are the most challenging, resulting in longer solving times and lower success rates among participants.

## 6.2 Real-world Performance of Our Web Application

**Setup.** To evaluate our web application’s real-world performance, we conducted a follow-up study with the same 10 participants. Unlike the previous study where participants solved pre-generated tasks from different sources, here each participant requested and solved 5 tasks in real-time through our web application. Participants selected their preferred contexts and made task requests until receiving a task. For each received task, we maintained the same instructions, time limits, and feedback questions as in the previous study.

**Results.** Our web application successfully synthesized tasks for 50/54 requests (92.6% coverage). Figure 9 provides insights into each participant’s session. Participants managed to solve on average 90.0% of the synthesized tasks, with an average solving time of 8.85 minutes. When analyzing their feedback, synthesized tasks showed on average high alignment with chosen theme (1.0) and programming concepts (0.95), while maintaining good comprehensibility (0.86).

## 7 Concluding Discussions

We introduced PYTASKSYN, a novel synthesis technique that leverages generative models as agents simulating different classroom roles to validate generated programming tasks. Each stage in our pipeline contributes uniquely to the validation process, collectively ensuring the creation of tasks that are of high-quality and meet diverse learning objectives. Through extensive expert evaluation and user studies, we demonstrated that our approach significantly improves the quality of generated programming tasks with minimal cost, while maintaining reasonable coverage across various themes and programming concepts.

Our work brings two important implications for leveraging generative AI for computing education. First, we demonstrate how to effectively automate the quality assessment of generated programming tasks. This paves the way to reducing educator workload when designing practice exercises that target specific learning objectives in programming education. Second, our results show that by breaking down the task synthesis process into different stages, we can leverage generative models that simulate different agents for specialized roles. This agent-based approach opens up new opportunities for utilizing generative models to simulate learning analytics to generate and validate educational content.

Next, we discuss some limitations of our work and possible ways of approaching them in the future. First, we adopt an accept/reject approach during our validation; future work could employ a framework for refining programming tasks by leveraging feedback from the validation stage. Second, we do not analyze whether the passing threshold of simulated students affects the synthesized task’s difficulty; it would be valuable to investigate it and see whether it aligns with difficulty assessed by experts. Third, we focus on Python programming; it would be interesting to explore whether the capabilities of generative models for validating tasks extend beyond Python. Fourth, we conducted studies with a relatively small number of participants; it would be important to conduct larger-scale studies with students and assess their learning outcomes.

## Acknowledgments and Disclosure of Funding

Funded/Co-funded by the European Union (ERC, TOPS, 101039090). Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union or the European Research Council. Neither the European Union nor the granting authority can be held responsible for them.

## References

- [1] Paul Denny, Sumit Gulwani, Neil T. Heffernan, Tanja Käser, Steven Moore, Anna N. Rafferty, and Adish Singla. Generative AI for Education (GAIED): Advances, Opportunities, and Challenges. *CoRR*, abs/2402.01580, 2024.
- [2] James Prather, Paul Denny, Juho Leinonen, Brett A. Becker, Ibrahim Albluwi, Michelle Craig, Hieke Keuning, Natalie Kiesler, Tobias Kohn, Andrew Luxton-Reilly, Stephen MacNeil, Andrew Petersen, Raymond Pettit, Brent N. Reeves, and Jaromír Savelka. The Robots Are Here: Navigating the Generative AI Revolution in Computing Education. In *Proceedings of the Working Group Reports of the Conference on Innovation and Technology in Computer Science Education (ITiCSE)*, 2023.
- [3] OpenAI. Hello GPT-4o. <https://openai.com/index/hello-gpt-4o/>, 2024.
- [4] GitHub. GitHub Copilot: Your AI Pair Programmer. <https://github.com/features/copilot>, 2022.
- [5] Hasan Abu-Rasheed, Christian Weber, and Madjid Fathi. Knowledge Graphs as Context Sources for LLM-Based Explanations of Learning Recommendations. In *Proceedings of the Global Engineering Education Conference (EDUCON)*, 2024.
- [6] Juho Leinonen, Arto Hellas, Sami Sarsa, Brent N. Reeves, Paul Denny, James Prather, and Brett A. Becker. Using Large Language Models to Enhance Programming Error Messages. In *Proceedings of the Technical Symposium on Computer Science Education (SIGCSE)*, 2023.
- [7] Tung Phung, Victor-Alexandru Padurean, Anjali Singh, Christopher Brooks, José Cambroner, Sumit Gulwani, Adish Singla, and Gustavo Soares. Automating Human Tutor-Style Programming Feedback: Leveraging GPT-4 Tutor Model for Hint Generation and GPT-3.5 Student Model for Hint Validation. In *Proceedings of the Learning Analytics and Knowledge Conference (LAK)*, 2024.
- [8] Hussein Mozannar, Gagan Bansal, Adam Fourney, and Eric Horvitz. Reading Between the Lines: Modeling User Behavior and Costs in AI-Assisted Programming. In *Proceedings of the CHI Conference on Human Factors in Computing Systems, (CHI)*, 2024.
- [9] Tung Phung, Victor-Alexandru Padurean, José Cambroner, Sumit Gulwani, Tobias Kohn, Rupak Majumdar, Adish Singla, and Gustavo Soares. Generative AI for Programming Education: Benchmarking ChatGPT, GPT-4, and Human Tutors. In *Proceedings of the Conference on International Computing Education Research (ICER) - Volume 2*, 2023.
- [10] Pragnya Sridhar, Aidan Doyle, Arav Agarwal, Christopher Bogart, Jaromír Savelka, and Majd Sakr. Harnessing LLMs in Curricular Design: Using GPT-4 to Support Authoring of Learning Objectives. In *AIED’23 Workshop on Empowering Education with LLMs - the Next-Gen Interface and Content Generation*, 2023.
- [11] Andre Del Carpio Gutierrez, Paul Denny, and Andrew Luxton-Reilly. Automating Personalized Parsons Problems with Customized Contexts and Concepts. In *Proceedings of the Conference on Innovation and Technology in Computer Science Education (ITiCSE)*, 2024.

- [12] Mollie Jordan et al. Need a Programming Exercise Generated in Your Native Language? ChatGPT’s Got Your Back: Automatic Generation of Non-English Programming Exercises Using OpenAI GPT-3.5. In *Proceedings of the Technical Symposium on Computer Science Education (SIGCSE)*, 2024.
- [13] Evanfiya Logacheva, Arto Hellas, James Prather, Sami Sarsa, and Juho Leinonen. Evaluating Contextually Personalized Programming Exercises Created with Generative AI. In *Proceedings of the Conference on International Computing Education Research (ICER)*, 2024.
- [14] Sami Sarsa, Paul Denny, Arto Hellas, and Juho Leinonen. Automatic Generation of Programming Exercises and Code Explanations Using Large Language Models. In *Proceedings of the Conference on International Computing Education Research (ICER)*, 2022.
- [15] Muhammad Fawad Akbar Khan, Max Ramsdell, Ha Nguyen, and Hamid Karimi. Human Evaluation of GPT for Scalable Python Programming Exercise Generation. In *International Conference on Data Science and Advanced Analytics (DSAA)*, pages 1–10. IEEE, 2024.
- [16] Andre Del Carpio Gutierrez, Paul Denny, and Andrew Luxton-Reilly. Evaluating Automatically Generated Contextualised Programming Exercises. In *Proceedings of the Technical Symposium on Computer Science Education (SIGCSE)*, 2024.
- [17] Jie Huang, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song, and Denny Zhou. Large Language Models Cannot Self-Correct Reasoning Yet. In *The Twelfth International Conference on Learning Representations, (ICLR)*, 2024.
- [18] OpenAI. GPT-4o mini: Advancing Cost-Efficient Intelligence. <https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/>, 2024.
- [19] Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chi-Min Chan, Heyang Yu, Yaxi Lu, Yi-Hsin Hung, Chen Qian, Yujia Qin, Xin Cong, Ruobing Xie, Zhiyuan Liu, Maosong Sun, and Jie Zhou. AgentVerse: Facilitating Multi-Agent Collaboration and Exploring Emergent Behaviors. In *International Conference on Learning Representations (ICLR)*, 2024.
- [20] Yilun Du, Shuang Li, Antonio Torralba, Joshua B. Tenenbaum, and Igor Mordatch. Improving Factuality and Reasoning in Language Models through Multiagent Debate. In *International Conference on Learning Representations (ICLR)*, 2024.
- [21] Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang, Xiaoyun Zhang, and Chi Wang. AutoGen: Enabling Next-Gen LLM Applications via Multi-Agent Conversation Framework. *CoRR*, abs/2308.08155, 2023.
- [22] Gati V. Aher, Rosa I. Arriaga, and Adam Tauman Kalai. Using Large Language Models to Simulate Multiple Humans and Replicate Human Subject Studies. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2023.
- [23] Xinyi Lu and Xu Wang. Generative Students: Using LLM-Simulated Student Profiles to Support Question Item Evaluation. In *Proceedings of the Conference on Learning @ Scale (L@S)*, 2024.
- [24] Tung Phung, José Cambronero, Sumit Gulwani, Tobias Kohn, Rupak Majumdar, Adish Singla, and Gustavo Soares. Generating High-Precision Feedback for Programming Syntax Errors using Large Language Models. In *Proceedings of the International Conference on Educational Data Mining (EDM)*, 2023.
- [25] Manh Hung Nguyen, Sebastian Tschiatschek, and Adish Singla. Large Language Models for In-Context Student Modeling: Synthesizing Student’s Behavior in Visual Programming from One-Shot Observation. In *Proceedings of the International Conference on Educational Data Mining (EDM)*, 2023.
- [26] Umair Z. Ahmed, Maria Christakis, Aleksandr Efremov, Nigel Fernandez, Ahana Ghosh, Abhik Roychoudhury, and Adish Singla. Synthesizing Tasks for Block-based Programming. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2020.
- [27] Victor-Alexandru Pădurean, Georgios Tzannetos, and Adish Singla. Neural Task Synthesis for Visual Programming. *Transactions on Machine Learning Research (TMLR)*, 2024.

- [28] Chao Wen, Ahana Ghosh, Jacqueline Staub, and Adish Singla. Task Synthesis for Elementary Visual Programming in XLogoOnline Environment. In *Proceeding of the International Conference on Artificial Intelligence in Education AIED*, 2024.
- [29] Darren Key, Wen-Ding Li, and Kevin Ellis. Toward Trustworthy Neural Program Synthesis. *CoRR*, abs/2210.00848, 2022.
- [30] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging LLM-as-a-Judge with MT-Bench and Chatbot Arena. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2023.
- [31] Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. Teaching Large Language Models to Self-Debug. In *International Conference on Learning Representations (ICLR)*, 2024.
- [32] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegraffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. Self-Refine: Iterative Refinement with Self-Feedback. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2023.
- [33] Zheyuan Zhang, Daniel Zhang-Li, Jifan Yu, Linlu Gong, Jinchang Zhou, Zhiyuan Liu, Lei Hou, and Juanzi Li. Simulating Classroom Education with LLM-Empowered Agents. *CoRR*, abs/2406.19226, 2024.
- [34] Julia M. Markel, Steven G. Opferman, James A. Landay, and Chris Piech. GPTeach: Interactive TA Training with GPT-based Students. In *Proceedings of the Conference on Learning @ Scale (L@S)*, 2023.
- [35] Hyounghwook Jin, Seonghee Lee, Hyungyu Shin, and Juho Kim. Teach AI How to Code: Using Large Language Models as Teachable Agents for Programming Education. In *Proceedings of the International Conference on Human Factors in Computing Systems (CHI)*, 2024.
- [36] Qianou Ma, Hua Shen, Kenneth Koedinger, and Sherry Tongshuang Wu. How to Teach Programming in the AI Era? Using LLMs as a Teachable Agent for Debugging. In *Proceeding of the International Conference on Artificial Intelligence in Education (AIED)*, 2024.
- [37] Robin Schmucker, Meng Xia, Amos Azaria, and Tom M. Mitchell. Ruffle&Riley: Towards the Automated Induction of Conversational Tutoring Systems. In *NeurIPS’23 Workshop on Generative AI for Education (GAIED)*, 2023.
- [38] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2022.
- [39] Jacob Cohen. A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20:37 – 46, 1960.