Large Language Models in Introductory Programming Education: ChatGPT's Performance and Implications for Assessments

Natalie Kiesler (D) Daniel Schiffner (D²)

Abstract: This paper investigates the performance of the Large Language Models (LLMs) ChatGPT-3.5 and GPT-4 in solving introductory programming tasks. Based on the performance, implications for didactic scenarios and assessment formats utilizing LLMs are derived. For the analysis, 72 Python tasks for novice programmers were selected from the free site CodingBat. Full task descriptions were used as input to the LLMs, while the generated replies were evaluated using CodingBat's unit tests. In addition, the general availability of textual explanations and program code was analyzed. The results show high scores of 94.4 to 95.8% correct responses and reliable availability of textual explanations and program code, which opens new ways to incorporate LLMs into programming education and assessment.

Keywords: Large language models; ChatGPT-3.5; GPT-4; Conversational Programming; Assessment

1 Introduction

The advent of Large Language Models (LLMs), such as OpenAI's ChatGPT, Codex, and GitHub's Copilot, affects the educational landscape at its core, as LLMs offer entirely new possibilities, but also challenges for educators, learners, and institutions. Even though LLMs have only appeared very recently to a broader audience, research has started to address their implications on computing education, particularly programming. The generative potential may be used by educators for the design of new programming tasks [Sa22], or for students to gather formative feedback [Ka23, Zh22]. At the same time, implications for programming pedagogy and assessments are being discussed [Be23, BK23, RTT23], as the low-threshold availability of LLMs raises new questions with regard to adequate task designs, students' contribution, plagiarism, and ethical conduct. Educators and institutions will soon need to reconsider the design of (formative) assessments. In this context, it is crucial to investigate the capabilities and limitations of LLMs for novice learners of programming, whose challenges have a well-documented history [SS86, Mc01, Lu18].

 $^{^{1}\,\}mathrm{DIPF}$ Leibniz Institute for Research and Information in Education, Germany, kiesler@dipf.de

² DIPF Leibniz Institute for Research and Information in Education, Germany, schiffner@dipf.de

The goals of this paper are (1) to investigate the potential of LLMs, such as Chat-GPT, for the generation of correct, executable program code for introductory programming tasks, and (2) to discuss didactic scenarios including assessments for the use of LLMs in introductory programming education.

2 Related Work on Large Language Models

Research on Large Language Models (LLMs) started to increase ever since OpenAI's ChatGPT was launched with free access in November 2022. Other LLMs of interest for the context of computing and programming education comprise OpenAI's Codex, and GitHub's Copilot³. Early papers on the performance of OpenAI's Codex conclude that the code generated by Codex outscores most students in CS1 [Fi22] and CS2 exercises [Fi23]. Zhang et al. [Zh22] found that Codex can help students to fix syntactic and semantic mistakes in their Python code. In a study with 69 programming novices, Kazemitabaar et al. [Ka23] explored the potential of Codex for solving programming tasks. Students with access to Codex had completed their tasks and increased their scores significantly, compared to those students without access to Codex. GitHub's Copilot is also capable of successfully solving introductory programming tasks with few additional prompts or adjustments [We23, PS22].

However, the code and feedback provided by LLMs along with usability issues imply that there is still room for improvements and further developments. For example, a study with 24 programmers concluded that the code generated by Copilot still contains errors [VZG22]. Even though test subjects preferred Copilot over the code completion plugin IntelliSense, longer code snippets were perceived as difficult to understand, edit and repair from the programmer's perspective. Similarly, Prather et al. [Pr23b] identified challenges of novices when using Copilot, some of which are due to Copilot's design. A study by Denny et al. [DKG23] focused on the engineering of prompts when using Copilot to mitigate its performance deficits and identify prompts leading to better feedback and results.

An exploratory interview study with five professional developers [BK23] summarizes opportunities and threats from an industry perspective, and discusses implications on software development education. Among them are scaffolding and fading of support based on Sweller's cognitive load theory [SAK11], a change of assessment formats, and a transitional period for novice learners before using LLMs. Rudolph et al. [RTT23] discuss ChatGPT's impact on traditional assessment formats in higher education. They recommend, e.g., fostering and assessing students' creative and critical thinking abilities, and letting students perform their competencies in class, or in authentic situations, while offering choices (if possible) to

³ available here: https://openai.com/blog/openai-codex and here: https://github.com/features/copilot

address students' interests. The conclusion is "to help students learn how to use AI tools judiciously and understand their benefits and limitations" [RTT23]. This is also the focus of a recent working group within the context of the Innovation and Technology in Computer Science Education conference [Pr23a].

3 Methodology

To address the current research desiderata on ChatGPT's performance in introductory programming tasks, and opportunities related to conversational programming in education and assessment, we define the following research question: How does a Large Language Model like ChatGPT perform when asked to solve introductory programming tasks?

To address this research question, all CodingBat tasks from Python code practice areas were used as input to GPT-3.5 (freely available), and GPT-4 (pay model) as the most recent versions of OpenAI's LLM. The CodingBat tasks are available within eight areas and contain simple, basic, and medium problems. Each area contains 6 to 12 exercises, resulting in a total of 72 tasks [Pa17]. The advantage of CodingBat tasks is that they are available online, and have been authored by a computing educator from Stanford University. Moreover, the automatic evaluation of input via unit tests upon execution allows for a straightforward evaluation of the LLM's performance. For every problem, the exact task description is presented to the LLM, (see [DKG23]). Then the generated output is characterized using a simple metric (e.g., textual explanation, or program code contained). If ChatGPT's response contained program code in Python, the code was used as input to CodingBat. Then the suggested program was executed, and the number of correct and incorrect test cases was recorded. In case ChatGPT could not provide a correct solution, additional prompts were explored as input and, once again, evaluated w.r.t. test case performance via CodingBat.

4 LLMs' Performance in Solving Introductory Programming Tasks

We evaluated both Chat-GPT-3.5 and GPT-4 and their performance in introductory programming tasks provided by CodingBat [Pa17]. Upon entering the full task description as an input, ChatGPT-3.5 immediately solved the task correctly in 69 out of 72 cases (95,8%), while GPT-4 solved 68 out of 72 tasks correctly (94,4%). Program code in Python was provided in all initial responses by both, ChatGPT-3.5 and GPT-4, even though the tasks do not explicitly ask for it. Textual explanations of the code were generated for 70 out of 72 tasks by both, ChatGPT-3.5 and GPT-4, and thus in 97,2% of all test cases. Moreover, the code often contains

```
python

def alarm_clock(day, vacation):
    # Check if on vacation
    if vacation:
        # If on vacation, check if weekend
        if day == 0 or day == 6:
            return 'off' # 'off' on weekends during vacation
    else:
        return '10:00' # '10:00' on weekdays during vacation
else:
        # If not on vacation, check if weekend
        if day == 0 or day == 6:
            return '10:00' # '10:00' on weekends
        else:
            return '7:00' # '7:00' on weekdays
```

Fig. 1: Code with comments for the alarm clock task generated by ChatGPT-3.5

explanatory comments (see Figure 1), or additional sample output is generated. Table 1 presents the performance results of both LLM versions in each of the eight task areas. Thus, the availability of textual explanations, program code, and the generation of fully correct solutions passing all unit tests are summarized for the number of tasks in each area.

Tab 1: Summary	of CPT 3.5's and	GPT-4's performance	in colving	CodingBat tacks
rab. 1: Summary	OF GET 1-9.0 S and	GP 1-4 S Demormance	III SOLVIII9	Coump Dat Jasks.

CodingBat task area	GPT-3.5 textual explana- tion	GPT- 3.5 pro- gram code	GPT-3.5 correct unit test results	GPT-4 textual explana- tion	GPT-4 pro- gram code	GPT-4 correct unit test results
Warmup1	11/12	12/12	12/12	12/12	12/12	12/12
Warmup2	9/9	9/9	9/9	9/9	9/9	9/9
String1	11/11	11/11	10/11	11/11	11/11	11/11
List1	12/12	12/12	12/12	11/12	12/12	12/12
Logic1	8/9	9/9	8/9	9/9	9/9	9/9
Logic2	7/7	7/7	6/7	6/7	7/7	6/7
String2	6/6	6/6	6/6	6/6	6/6	4/6
List2	6/6	6/6	6/6	6/6	6/6	5/6

The three errors made by Chat-GPT-3.5 were due to an additionally required method (task $make_tags$), and ambiguity or a lack of clarity in the task description (task $near_ten$, and $make_bricks$). The four errors in the responses generated by GPT-4 were due to similar reasons: ambiguity in the task description (task $make_chocolate$, and $count_hi$). The reason why two other responses caused the feedback "Bad code" was that GPT-4 imported libraries, which is a general constraint among CodingBat tasks. These errors were observed for $count_code$ and $centered_average$.

To improve the generated solutions, the following prompts were used: "Please generate compilable Python code" upon compile problems (task <code>make_tags</code>), and "The code fails the following test cases: <test cases>" for failed test cases (e.g., for <code>near_ten</code>). These prompts immediately resulted in an improved, correct response including an explanation. However, when asked to correct the response to the <code>make_bricks</code> task, ChatGPT-3.5 was reluctant to change its output. The prompt including the incorrect test cases had to be repeated three times, before resulting in a fully correct solution. This seemingly overconfidence was not observed with GPT-4. The latter immediately reached a fully correct solution after using the prompt with the failed test cases (<code>make_chocolate</code>, and <code>count_hi</code>), or the request to "Please solve the task without an import." (<code>count_code</code> and <code>centered_average</code>) once.

Despite high rates of successful task completion, the performance has to be discussed w.r.t. to several aspects. First of all, the selected tasks are straightforward and clear in most cases, as they were developed by an experienced educator. They contain little ambiguity and provide exemplary input and output. It was thus a successful strategy to use the full task description as input to the LLM. Nonetheless, students should be aware that ambiguity in the task will likely cause incorrect responses, as ChatGPT may offer a solution to a different problem. Now this does not mean that teachers should develop ambiguous tasks. On the contrary, clear assignments are still important so that students understand the task.

Second, students need to adhere to general task constraints. For instance, many educators do not allow the use of libraries, or they provide certain function signatures for novice learners. The same is true for CodingBat tasks. Hence, students must be aware of such more or less explicit constraints, and provide them as additional information to the LLM, if they want to receive correct answers. Moreover, students need to know that ChatGPT can be overly confident, as it may not immediately change its proposed solution. We observed this phenomenon while using ChatGPT-3.5, but not with GPT-4. Providing failed test cases as a follow-up input seems to be an adequate strategy to improve the output, but fully trusting an LLM is not (yet) an option.

A limitation of this work is due to the random nature of ChatGPT's responses. Thus, answers are arbitrary, and we doubt that the very same answers can be replicated by other researchers. Moreover, model solutions to all CodingBat tasks are available in GitHub repositories, so chances are ChatGPT was trained on such data. In addition, CodingBat only offers very small programs, with real novices as a target group, meaning that common second-semester tasks are very different, and so might be the LLM performance. Thus, continuous research on the evolving capabilities of LLMs is required to evaluate their implications on education for higher semesters.

5 Discussion of Implications on Didactic Scenarios and Assessments

After having discussed ChatGPT's current performance (June 2023) in solving programming tasks, we now focus on the implications of these results on didactic scenarios and assessment formats in introductory programming courses. We, therefore, propose some exemplary settings for the inclusion of LLMs in learning activities and formative assessments.

As a guiding principle, we evaluated the results and didactic scenarios using the following rule: Assume that the response is invalid or contains errors and that the LLM may be overconfident and hallucinating. Considering the good performance for the analyzed tasks and having this rule in mind, it seems quite natural that current LLMs can be used by students for individual practice and self-assessment. The option to discuss a solution can, for example, provide valuable information to learners. This might be comparable to a peer-review session or direct input/feedback from an educator. Even though responses may not be perfect, LLMs can provide useful textual explanations and code suggestions students currently gather from other sources, e.g., Stack Overflow.

As our examples showed high success rates for the given tasks, novice programmers could also use the generated code to compare it to their own solution. A simple request to compare solutions provides the rationale for why one solution may be preferable over another one. Even simple tasks can have multiple solutions, but it is not obvious to novice programmers when to choose one over the other. Interesting results from such an exercise may also be discussed further in class with additional elaboration by the facilitator.

Another concept is to let students generate multiple solutions by LLMs, such as ChatGPT on purpose, and to discuss them as part of a peer-group exercise w.r.t. advantages and disadvantages. GPT-3.5 and GPT-4 already generate different solutions in many cases, and the same is true upon regeneration of the reply. While they solve the given task, the discussion among peers gives students the option to understand the underlying problem class more thoroughly. In our experiments, GPT-4 created solutions that were more sophisticated and even more complex than necessary, for example, by including additional conditions (e.g., $same_first_last$). In other cases, GPT-3.5's responses were unnecessarily complex (e.g., sum_3). Such discussions can potentially be very fruitful and help understand several problem-solving approaches for one problem.

In another scenario, students' discussions with ChatGPT can be used for the evaluation of learning processes and students' mental models. Currently, we are limited to either a static/dynamic code analysis that builds upon a test-driven development approach. In contrast to that, strategies for problem-solving are rarely assessed.

With an LLM like ChatGPT and the option to share such a discussion via a link, a more individualized approach is available, helping teachers to better understand issues novice programmers are facing. LLMs like ChatGPT thus enable a new form of reflecting on problem-solving approaches to programming tasks. The discussion with such a tool can provide insights into the logic used to solve a given task. It allows educators to see that and distinguish between algorithmic and coding issues, which are hard to detect with program code alone.

Another formative assessment method may involve reflection exercises with a focus on critically discussing tasks and various program solutions with ChatGPT. This way, students can learn more about program code (i.e., develop knowledge about knowledge [AK01]), while critically approaching LLMs and their generated solutions. Conversations with LLMs may therefore be used as an innovative assessment method for meta-cognitive competencies, that is the systematic approach towards solving (similar) problems [Ki20b]. LLMs can further facilitate new ways to assess lower levels of cognitive complexity as it becomes easier to represent tasks aiming at the analysis and evaluation of code, which, in turn, can contribute to the development and transfer of problem-solving strategies. The same may apply to the understanding of seemingly simple programming concepts, or procedural knowledge (see [Ki22, Ki20a, Ki20c] for the classification of programming competencies). A prerequisite for such an assessment is, of course, the identification of observable and reliable indicators for such a measurement, along with an introduction of learners to the potential and limitations of conversing with LLMs (see also [BK23, RTT23]).

6 Conclusion and Future Work

In this study, we investigated the potential of large language models in introductory programming education and assessments. The research question addresses the performance of LLMs when asked to solve introductory programming tasks. To answer it, we utilized existing coding tasks, i.e., CodingBat, as input to ChatGPT-3.5 and GPT-4 to produce program code, which was then evaluated by CodingBat's unit tests. The results show a high success rate, ranging between 94.4 and 95.8%, which allowed us to further discuss didactic scenarios including assessments. Several scenarios, which until recently have been limited to having an expert (educator or tutor) at a learner's side, seem more realistic now – even for large courses. LLMs further enable the assessment of different cognitive process dimensions, such as understanding programming concepts, or reading and analyzing code, which are currently hard to implement, and rarely addressed in higher education programming courses. Overall, the availability of LLMs may be a blessing for education if used with caution and guidance. With this study, we showed that an application in the context of introductory programming tasks can be reasonable. We expect these tools to help novice programmers to better understand problems and concepts, and that we can overcome some of the current limitations in programming education and assessments.

Options for future work are manifold, as the full potential and challenges of LLMs have not yet been evaluated. Several follow-up studies are currently work-in-progress. One of them is concerned with the evaluation of feedback types offered by LLMs [KLK23], and another one aims to investigate ChatGPT's performance in actual exam tasks of an introductory programming course. The development of benchmarks for certain tasks is also considered future work [Pr23a], as this could help educators evaluate the adequacy of assessments and didactic scenarios.

Bibliography

- [AK01] Anderson, Lorin W.; Krathwohl, David R.: A Taxonomy for Learning, Teaching, and Assessing. A Revision of Bloom's Taxonomy of Educational Objectives. Addison Wesley Longman, New York, 2001.
- [Be23] Becker, Brett A; Denny, Paul; Finnie-Ansley, James; Luxton-Reilly, Andrew; Prather, James; Santos, Eddie Antonio: Programming Is Hard – Or at Least It Used to Be: Educational Opportunities And Challenges of AI Code Generation. ACM, 2023.
- [BK23] Bull, Christopher; Kharrufa, Ahmed: , Generative AI Assistants in Software Development Education, 2023.
- [DKG23] Denny, Paul; Kumar, Viraj; Giacaman, Nasser: 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. SIGCSE 2023, ACM, USA, p. 1136-1142, 2023.
- [Fi22] Finnie-Ansley, James; Denny, Paul; Becker, Brett A; Luxton-Reilly, Andrew; Prather, James: The robots are coming: Exploring the implications of openai codex on introductory programming. In: Australasian Computing Education Conference. pp. 10–19, 2022.
- [Fi23] Finnie-Ansley, James; Denny, Paul; Luxton-Reilly, Andrew; Santos, Eddie Antonio; Prather, James; Becker, Brett A.: My AI Wants to Know if This Will Be on the Exam: Testing OpenAI's Codex on CS2 Programming Exercises. ACM, pp. 97–104, 1 2023.
- [Ka23] Kazemitabaar, Majeed; Chow, Justin; Ma, Carl Ka To; Ericson, Barbara J; Weintrop, David; Grossman, Tovi: 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. pp. 1–23, 2023.
- [Ki20a] Kiesler, Natalie: On Programming Competence and Its Classification. In: Koli Calling '20: Proceedings of the 20th Koli Calling International Conference on Computing Education Research. Koli Calling '20, Association for Computing Machinery, New York, NY, USA, 2020.

- [Ki20b] Kiesler, Natalie: Towards a Competence Model for the Novice Programmer Using Bloom's Revised Taxonomy - An Empirical Approach. In: Proceedings of the 2020 ACM Conference on Innovation and Technology in Computer Science Education. ITiCSE '20, Association for Computing Machinery, New York, NY, USA, p. 459–465, 2020.
- [Ki20c] Kiesler, Natalie: Zur Modellierung und Klassifizierung von Kompetenzen in der grundlegenden Programmierausbildung anhand der Anderson Krathwohl Taxonomie. In: CoRR abs/2006.16922. arXiv: 2006.16922. url: https://doi.org/10.48550/arXiv.2006.16922. 2020.
- [Ki22] Kiesler, Natalie: Kompetenzförderung in der Programmierausbildung durch Modellierung von Kompetenzen und informativem Feedback. Dissertation, Johann Wolfgang Goethe-Universität, Frankfurt am Main, Januar 2022. Fachbereich Informatik und Mathematik.
- [KLK23] Kiesler, Natalie; Lohr, Dominic; Keuning, Hieke: Exploring the Potential of Large Language Models to Generate Formative Programming Feedback. In: 2023 IEEE Frontiers in Education Conference (FIE), College Station, TX, USA. pp. 1–5, 2023.
- [Lu18] Luxton-Reilly, Andrew; Simon; Albluwi, Ibrahim; Becker, Brett A.; Giannakos, Michail; Kumar, Amruth N.; Ott, Linda; Paterson, James; Scott, Michael James; Sheard, Judy; Szabo, Claudia: Introductory Programming: A Systematic Literature Review. In: Proceedings of the 23rd Annual ACM Conference on Innovation and Technology in Computer Science Education. ACM, New York, pp. 55–106, 2018.
- [Mc01] McCracken, Michael; Almstrum, Vicki; Diaz, Danny; Guzdial, Mark; Hagan, Dianne; Kolikant, Yifat Ben-David; Laxer, Cary; Thomas, Lynda; Utting, Ian; Wilusz, Tadeusz: A Multi-National, Multi-Institutional Study of Assessment of Programming Skills of First-Year CS Students. In: Working Group Reports from ITiCSE on Innovation and Technology in Computer Science Education. ACM, New York, p. 125–180, 2001.
- [Pa17] Parlante, Nick: CodingBat Python. https://codingbat.com/python, 2017.
- [Pr23a] Prather, James; Denny, Paul; Leinonen, Juho; Becker, Brett A.; Albluwi, Ibrahim; Caspersen, Michael E.; Craig, Michelle; Keuning, Hieke; Kiesler, Natalie; Kohn, Tobias; Luxton-Reilly, Andrew; MacNeil, Stephen; Petersen, Andrew; Pettit, Raymond; Reeves, Brent N.; Savelka, Jaromir: Transformed by Transformers: Navigating the AI Coding Revolution for Computing Education: An ITiCSE Working Group Conducted by Humans. In: Proceedings of the 2023 Conference on Innovation and Technology in Computer Science Education V. 2. ITiCSE 2023, Association for Computing Machinery, New York, p. 561–562, 2023.
- [Pr23b] Prather, James; Reeves, Brent N; Denny, Paul; Becker, Brett A; Leinonen, Juho; Luxton-Reilly, Andrew; Powell, Garrett; Finnie-Ansley, James; Santos, Eddie Antonio: Ït's Weird That it Knows What I Want": Usability and Interactions with Copilot for Novice Programmers. 2023.
- [PS22] Puryear, Ben; Sprint, Gina: Github Copilot in the Classroom: Learning to Code with AI Assistance. J. Comput. Sci. Coll., 38(1):37–47, nov 2022.

- [RTT23] Rudolph, Jürgen; Tan, Samson; Tan, Shannon: ChatGPT: Bullshit spewer or the end of traditional assessments in higher education? Journal of Applied Learning and Teaching, 6(1), 2023.
- [Sa22] Sarsa, Sami; Denny, Paul; Hellas, Arto; Leinonen, Juho: Automatic Generation of Programming Exercises and Code Explanations Using Large Language Models. In: Proceedings of the 2022 ACM Conference on International Computing Education Research. ACM, New York, p. 27–43, 2022.
- [SAK11] Sweller, John; Ayres, Paul; Kalyuga, Slava: The Worked Example and Problem Completion Effects. In: Cognitive Load Theory. Springer, New York, pp. 99– 109, 2011.
- [SS86] Spohrer, James C.; Soloway, Elliot: Novice mistakes: Are the folk wisdoms correct? Communications of the ACM, 29(7):624-632, 1986.
- [VZG22] Vaithilingam, Priyan; Zhang, Tianyi; Glassman, Elena L: Expectation vs. experience: Evaluating the usability of code generation tools powered by large language models. pp. 1–7, 2022.
- [We23] Wermelinger, Michel: Using GitHub Copilot to Solve Simple Programming Problems. In: Proceedings of the 54th ACM Technical Symposium on Computer Science Education V. 1. SIGCSE 2023, ACM, USA, p. 172–178, 2023.
- [Zh22] Zhang, Jialu; Cambronero, José; Gulwani, Sumit; Le, Vu; Piskac, Ruzica; Soares, Gustavo; Verbruggen, Gust: Repairing Bugs in Python Assignments Using Large Language Models. 2022.