

Human-AI Co-Creation of Worked Examples for Programming Classes

Mohammad Hassany¹, Peter Brusilovsky¹, Jiaze Ke², Kamil Akhuseyinoglu¹ and Arun Balajee Lekshmi Narayanan¹

¹University of Pittsburgh, Pittsburgh, PA, 15260

²Carnegie Mellon University, Pittsburgh, PA, 15213

Abstract

Worked examples (solutions to typical programming problems presented as a source code in a certain language and are used to explain the topics from a programming class) are among the most popular types of learning content in programming classes. Most approaches and tools for presenting these examples to students are based on line-by-line explanations of the example code. However, instructors rarely have time to provide line-by-line explanations for a large number of examples typically used in a programming class. In this paper, we explore and assess a human-AI collaboration approach to authoring worked examples for Java programming. We introduce an authoring system for creating Java worked examples that generates a starting version of code explanations and presents it to the instructor to edit if necessary. We also present a study that assesses the quality of explanations created with this approach.

Keywords

Code Examples, Authoring Tool, Human-AI Collaboration

1. Introduction

Program code examples play a crucial role in learning how to program [1]. Instructors use examples extensively to demonstrate the semantics of the programming language being taught and to highlight the fundamental coding patterns. Programming textbooks also pay a lot of attention to examples, with a considerable textbook space allocated to program examples and associated comments [2, 3]. A typical worked example presents a code for solving a specific programming problem and explains the role and function of code lines or code chunks. In textbooks, these explanations are usually presented as comments in the code or as explanations on the margins. While informative, this approach focused on passive learning, which is known for its low efficiency. Recognizing this problem, several research teams developed learning tools that offered more interactive and engaging ways to learn from examples [4, 5, 6, 7, 8].

The example-focused learning tools demonstrated their effectiveness in classroom studies, but their use by programming instructors is still limited due to the insufficient number of worked

Joint Proceedings of the ACM IUI Workshops 2024, March 18-21, 2024, Greenville, South Carolina, USA

✉ moh70@pitt.edu (M. Hassany); peterb@pitt.edu (P. Brusilovsky); jiazek@andrew.cmu.edu (J. Ke);

kaa108@pitt.edu (K. Akhuseyinoglu); arl122@pitt.edu (A. B. L. Narayanan)


🌐 <https://github.com/mhassany-pitt/> (M. Hassany); <https://sites.pitt.edu/~peterb/> (P. Brusilovsky)

🆔 0009-0004-8893-8454 (M. Hassany); 0000-0002-1902-1464 (P. Brusilovsky); 0009-0003-3122-2298 (J. Ke);

0000-0002-7761-9755 (K. Akhuseyinoglu); 0000-0002-7735-5008 (A. B. L. Narayanan)



© 2024 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

 CEUR Workshop Proceedings (CEUR-WS.org)

examples offered by these tools. Although the authors of these tools usually provide a good set of worked examples that can be presented through their tools, many instructors prefer to use their own favorite code examples. The instructors are usually happy to broadly share the code of examples they created (usually providing it on the course web page), but they rarely have time or patience to augment examples with explanations and add their examples to an example-focused interactive system. Indeed, producing a single explained example could take 30 minutes or more, since it requires typing an explanation for each code line [4, 8] or creating a screencast in a specific format [5, 7].

This issue has been recognized by several research teams that have offered several ways to address the lack of content. Among the approaches explored are learner-sourcing, that is, engaging students in creating and reviewing explanations for instructor-provided code [9] and automatic extraction of information content from available sources, such as lecture recordings [6]. In this paper, we present an alternative approach to address the lack of worked examples based on human-AI collaboration. With this approach, the instructor provides the code of one of their favorite examples along with the statement of the programming problem it is solving. The AI engine based on large language models (LLM) examines the code and generates explanations for each code line. The explanations could be reviewed and edited by the instructor. To support and explore this authoring approach, we created an authoring system, which radically decreases the time to create a new interactive worked example. The examples created by the system could be uploaded to an example-exploration system such as WebEx [4] or PCEX [8] or exported in a reusable format. To assess the quality of the resulting examples, we performed a user study in which TAs and students compared code explanations created by experts through a traditional process with examples created by AI to contribute to human-AI collaborative process.

The remainder of the paper is structured as following. We start by reviewing related work, introduce the example authoring system that implements the proposed collaborative approach, and explain how specific design decisions were made through several rounds of internal evaluation. Next, we explain the design of our user study and review its results. We conclude with a summary of the work and plans for future research.

2. Related Work

2.1. Worked Examples in Programming

Code examples are important pedagogical tools for learning programming. Not surprisingly, considerable efforts have been devoted to the development of learning materials and tools to support students in studying code examples. For many years, the state-of-the-art approach for presenting worked code examples in online tools was simply code text with comments [1, 10, 11]. More recently, this approach has been enhanced with multimedia by adding audio narrations to explain the code [12] or by showing video fragments of code screencasts with the instructor's narration being heard while watching code in slides or an editor window [5, 6]. Both ways, however, support *passive* learning, which is the least efficient approach from the

prospect of the ICAP framework [13]¹

An attempt to make learning from program construction examples *active* was made in the WebEx system, which allowed students to interactively explore instructor-provided line-by-line comments for program examples via a web-based interface [4]. More recently, several projects [6, 7, 8] augmented examples with simple problems and other constructive activities to elevate the example study process to the *interactive and constructive* levels of the ICAP framework, known as the most pedagogically efficient.

A good example of a modern interactive tool for studying code examples is the PCEX system [8]. PCEX (Program Construction EXamples) was created in the context of an NSF Infrastructure project (<https://csssplice.org>) with a focus on broad reuse and has been used by several universities in the US and Europe in the context of Java, Python, and SQL courses. PCEX interface (Figure 1) provides interactive access to traditionally organized worked examples, i.e., code lines augmented with instructor’s explanations. Separating explanations (Figure 1-3) from the code (Figure 1-2), allows students to selectively study explanations for code lines they want. Explanations are provided on several levels of detail, so more details could be requested if the brief explanation is not sufficient (Figure 1-3).

Since line-by-line multi-level example explanations offered by PCEX is currently the most detailed approach for explaining worked examples, we selected the code example structure implemented by PCEX as the target model for our authoring tool presented in this paper. The tool produces code augmented with line-by-line explanations on several levels of detail. The resulting example could be directly uploaded to PCEX or exported in a system-independent format to be uploaded to other example exploration systems like WebEx [4].

2.2. Use of LLMs for Code Explanations

Several research teams explored the use of LLM for code explanations using GPT-3 [14, 15, 16], GPT-3.5 [15, 17, 18], GPT-4 [17], OpenAI Codex [19, 20, 15], and GitHub Copilot [18]. LLMs were used to generate explanations at different levels of abstraction (line-by-line, step-by-step, and high-level summary). Sarsa et al. [19] observed that ChatGPT can generate better explanations at low-level (lines). Explanations and summaries generated by these LLMs were mostly evaluated by authors [19], students [15, 16], and tool users [18]. Sarsa et al. [19] reported a high correct ratio for generated explanations with minor mistakes that can be resolved by the instructor or teaching assistant. Students rated LLM-generated explanations as being useful, easier, and more accurate than learner-sourced explanations [16].

Since prompts directly influence the LLM’s performance, several studies focused on exploring different prompting strategies [21, 22]. Tian et al [20] reported that a verbose prompt will limit the LLM’s ability to utilize its knowledge [20]. Iterative prompts are proven to perform well [14]. Zamfirescu-Pereira et al. [14] observed that non-experts have misconceptions about LLMs and struggle to come up with a well-formed prompt. Researchers believe that LLMs can be beneficial in environments where humans and AI can work together, where the human can perform the expert evaluation and tune the responses generated by the AI while the AI performs the time-consuming manual tasks [22].

¹The ICAP framework differentiates four modes of engagement, behaviorially exhibited by learners: *passive*, *active*, *constructive and interactive*.

Example: STUDENT: PointTester

Construct a class that represents a point in the Euclidean plane. The class should contain data that represents the point's integer coordinates (x,y). The class should also include getter and setter methods for accessing and changing the point's coordinates and a method to translate the point, i.e., shift the point's location by the specified amount.

The class PointTester1 instantiates an object from this class, sets the (x,y) coordinates of the point, and translates the point by the specified amount.

```

1 public class PointTester {
2     public static void main(String[] args) {
3         Point point = new Point();
4         point.setX(7);
5         point.setY(2);
6         point.translate(11, 6);
7         System.out.println("The point's coordinates: (" + point.getX() + ",
" + point.getY() + ")");
8     }
9 }
10 class Point {
11     private int x;
12     private int y;
13     public void translate(int dx, int dv) {

```

Challenge Me!

This line translates the point's location by shifting the x-coordinate by 11 and the y-coordinate by 6.

PREVIOUS NEXT

PREVIOUS ADDITIONAL DETAILS

Figure 1: Studying a code example in the PCEX system: 1) title and program description, 2) program source code with lines annotated with explanations, 3) explanations for the highlighted line, 4) link to a “challenge” - a small problem related to the example.

3. The Feasibility Studies

To assess the feasibility of Human-AI co-creation of worked examples, we performed three rounds of preliminary studies. The purpose of these studies was to develop an approach for producing LLM code explanations of reasonable quality, compare the explanations produced by LLMs with the explanations produced by humans, and assess whether the LLM explanations are considered satisfactory by instructors and students.

In the first study [23] guided by earlier work on LLM code explanations reviewed above, we explored a range of prompts and performed an evaluation of the quality of explanations generated by the prompts to select the best-performing prompt for the next rounds of our work.

In the second study [24], we used a dataset of explanations produced by two experts and 60 students for the same four Java code examples with 33 explainable lines to compare ChatGPT explanations with explanations produced by experts and students using several formal metrics. To make this comparison, we generated ChatGPT explanations using our selected prompt for the 33 explainable lines four times, using temperature 0 once and temperature 1 three times. To calculate all comparison metrics, we merged all line explanations generated by each source (i.e, each expert, each student, and each round of ChatGPT generation) into a single source document. As the data shows (Table 1), the explanations produced by ChatGPT have comparable length (measured by the number of tokens) and lexical density with the explanations produced by experts, while the explanations produced by students were more than twice as short and more lexically dense than the explanations produced by the other two sources. Surprisingly

(given the length difference) the readability of explanations produced by experts is very similar to the readability of student explanations, while ChatGPT explanations are much less readable. Expert explanations are also much more similar than ChatGPT explanations to the explanations produced by students (Table 2). This data could be partially explained by the considerably larger vocabulary used by ChatGPT even in comparison to experts.

Source	N	Vocabulary	Lexical Density	# of Tokens	GF	FRE	FK
Experts	2	209.0	0.48	690.0	8.46	78.45	6.18
ChatGPT*	4	238.0	0.49	769.5	11.09	69.64	7.83
Students	60	116.5	0.54	249.5	8.02	80.48	5.62

Table 1

Median lexical and readability metrics for different sources of explanations (FRE = Flesch-Reading Ease, FK = Flesch-Kincaid, GF = Gunning Fog). *refers to the prompt selected in the first study.

Reference	Source	chrF	METEOR	USE	BERTScore
Expert	Student	0.33	0.144	0.33	0.63
ChatGPT	Student	0.18	0.151	0.255	0.458
Expert	ChatGPT	0.32	0.28	0.48	0.712

Table 2

Assessing lexical and semantic alignment (larger is better) between sources of explanations.

In the third study [23], we conducted a comparative evaluation of explanations produced by experts and ChatGPT from the point of view of human users. We used two types of human users: authors (instructors and TAs) who are expected to use ChatGPT-generated explanations as the starting point in the co-creation process, and students who are the target users of the co-created product. Explanations were compared in pairs, each explanation in a pair has to be judged by completeness, and the best explanation in the pair has to be selected. A pair included an expert and a ChatGPT explanation, and the judges were not aware of which source produced each explanation. The study results indicated strong preferences for ChatGPT in both groups of judges (Table 3). In general, ChatGPT explanations were rated as more complete and judged to be better in the majority of cases. However, it was not a clear win. In a substantial number of cases (15.05% for students and 27.41% for authors), expert explanations were selected as the best option in a pair.

Taking the results of these two studies together, we could conclude that producing explanations for code examples is a promising application area for Human-AI co-creation. On the one hand, the LLM-generated explanations are lagging behind expert explanations in several aspects. ChatGPT explanations have higher reading difficulty than expert explanations, and they are further away from the students' own explanations, as measured by most similarity metrics. The vocabulary data hints that ChatGPT tends to use terms, which might not be easy for the students to understand, while experts have experience in phrasing their explanations closer to the students' active vocabulary. On the other hand, the explanations produced by ChatGPT were generally rated higher than the expert explanations by both instructors and students. These data hint that presenting ChatGPT explanations directly to students might not be a perfect solution, but they can serve as an excellent starting point for instructors in shaping

Source	Judged by	Not complete	Complete	Very complete	“This source is better”
ChatGPT	Students	0.00%	13.33%	86.67%	51.11%
ChatGPT	Authors	1.48%	32.59%	65.93%	58.15%
Experts	Students	2.22%	55.56%	42.22%	16.05%*
Experts	Authors	14.07%	57.78%	28.15%	27.41%*

Table 3

Assessment of explanations generated by ChatGPT and experts by students and authors. For convenience, we do not count the cases in which the explanations in a pair were judged equally good.

their own explanations. Following that, we decided to structure the Human-AI collaboration in creating working examples as follows. Instructors have the ultimate control over producing explanations. Depending on the context (such as example complexity), they can either choose to explain example lines themselves or request AI (LLM) help in producing explanations for specific lines. In the latter case, LLM generates the initial line explanations leaving it to the instructor to accept or reject it and, if accepted, to further edit the explanation text to satisfaction. The Human-AI co-creation interface presented in the next section is based on this model of collaboration.

4. The Human-AI Co-Creation Interface Design

On the basis of our feasibility studies, we developed a Worked Example Authoring Tool (WEAT). WEAT enables instructors to create worked code examples for PCEX system, [8] through the human-AI co-creation interface. In this co-creation process, the main task of a human author is to provide the code of the example and the statement of the problem that the code solves. The main task of ChatGPT is to generate the bulk of code line explanations on several levels of detail. As an option, a human author could edit and refine the text produced by ChatGPT to adapt it to the class goals and target students. As in any productive collaboration, each side does what it is best suited to do, leaving the rest to the partner.

In the main part of the WEAT interface, the problem (Figure 2-1) and the code (Figure 2-2) have to be provided by the instructor, while the explanations for each line (Figure 2-3) can be created by the instructor or generated by ChatGPT. The generated explanations could be further edited by the instructor. While we expect that co-creation of code explanations will be the preferred way to use WEAT, the system supports the whole range of options from using AI explanations without human editing to creating the whole example from scratch, without the help of AI. Authors who want to start by creating explanations themselves could simply select a code line to explain (Figure 2-2) and add one or more explanation fragments to this line (Figure 2-3). The order of the fragments is important: the first fragment is displayed in PCEX when the line is clicked, while the remaining fragments can be accessed by clicking the “Additional Details” button (Figure 1-3).

To generate ChatGPT explanations for the provided example code and problem description, the author has to click the “Generate Explanations” button to open the ChatGPT dialog (Figure 3). In this dialog, the explanations could be generated by clicking “Generate” button and added to

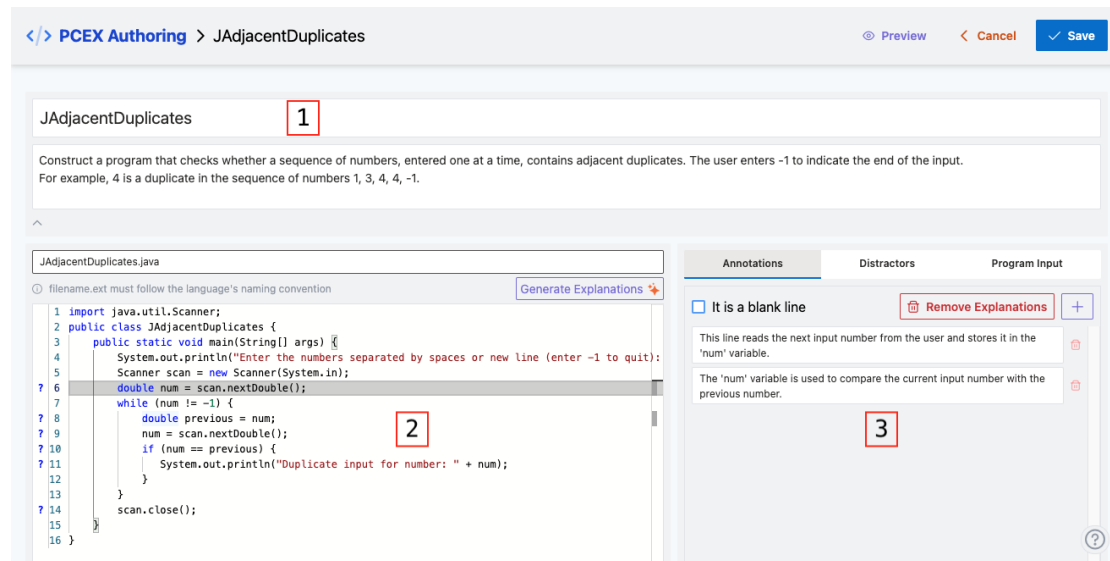


Figure 2: WEAT Authoring, 1) program title and description, 2) program source code (lines with explanations are marked with a blue question mark next to the line number), 3) explanations for the selected line (the line with gray background - line 6 in the screenshot).

the example by clicking “Use Explanations” button. Experienced authors have the opportunity to tune the default prompt before generating explanations and review the generated explanations before using them. Reviewing the generated explanations can be done line by line: selecting one of the explained lines (marked by “?”) in the code box (Figure 3-3) will display all generated explanations for this line in the explanation box (Figure 3-4). The explanation could be accepted or rejected by clicking the checkbox next to the “Include this line” prompt.

To support the review at the finer grain level, WEAT divides the explanations into fragments that can be independently accepted or rejected by clicking the small green check mark icon next to the fragment (Figure 3-4a). The author can also click on the small gray thumb-up icon (Figure 3-4b) to provide positive feedback on the explanation fragment. Once the “Use Explanations” button is clicked, all accepted explanation fragments are added to the corresponding example lines and can be further edited in the main interface (Figure 2).

5. Evaluation

To assess how well WEAT supports co-creation of worked examples, we engaged five instructors (A1-A5) teaching Java or Python classes and asked them to create one or more worked examples for PCEX from real examples they use in their classes. To explain the tool to the instructor, we provided a video tutorial and integrated textual help into WEAT. Their interactions and usage of the tool were recorded through logs and used for the analysis presented below.

The instructors used the tool to create 12 examples in total (Table 4). The ChatGPT dialog was used 21 times, and in 13 cases (A1=6, A2=2, A3=3, A4=1, and A5=1), instructors added generated explanations to the example by clicking the “Use Explanations” button. As discovered

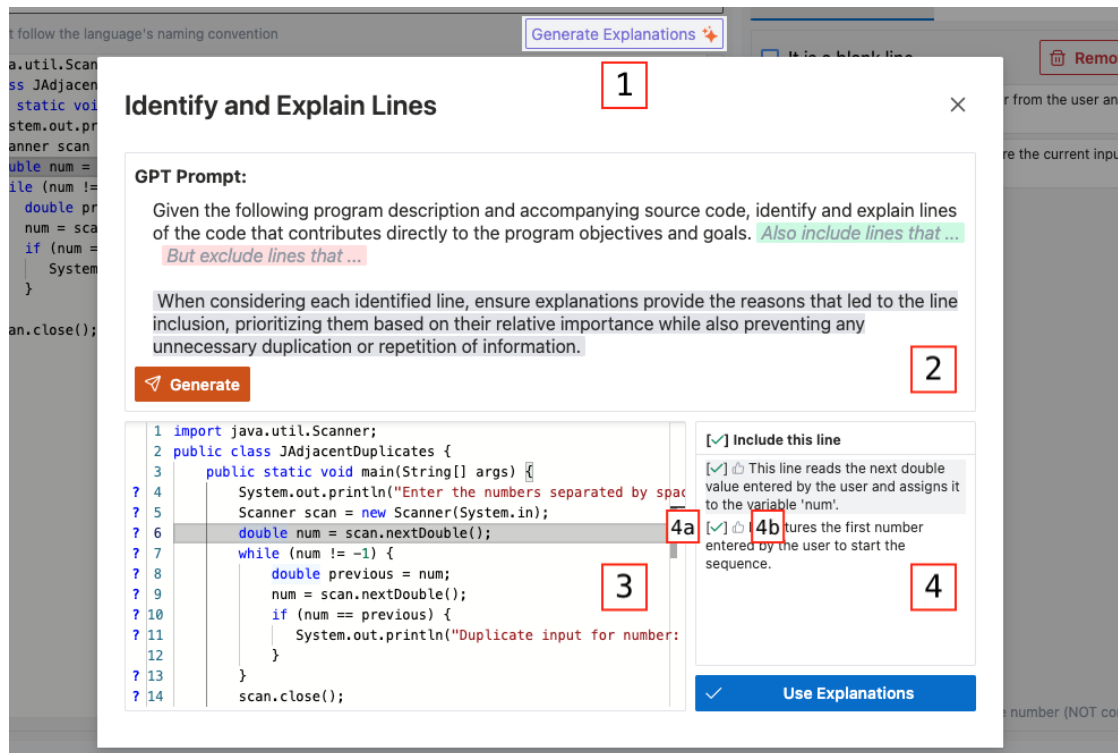


Figure 3: Human-AI Collaborative Worked Example Authoring, 1) “Generate Explanations” button, 2) default prompt (author can tune the prompt - optional), 3) program source preview, 4) generated explanations for the selected line.

from an interview with instructors, in several cases they closed and reopened the ChatGPT dialog to access the main interface blocked by the dialog. Analyzing the interaction logs, we observed this has been done at least 5 times (3 times with the close-reopen interval of 5 seconds and 2 with 12 seconds interval) leaving only 16 cases where explanations had a chance to be examined. In total, 269 explanation fragments were generated for 119 lines of code with an average of 2.26 fragments per line. In 13 cases where ChatGPT explanations were added to the example by instructors, ChatGPT generated 237 explanations for 99 lines of code (Table 4). We found no cases in which the entire set of explanations generated for the line was excluded by the instructors in its entirety, and among the 237 generated fragments, only 24 (10.12%) 237 were excluded. The interview revealed that in some cases the generated fragments were rejected not because they were unsatisfactory, but because they were incorrect (Figure 4). On the other hand, instructors liked 15 (6.32%) explanations.

After adding explanations to the example, instructors still didn’t remove the explanations for any line entirely, but removed 23 (9.7%) ChatGPT generated explanation fragments. Instructor A5 reported that he removed several fragments when merging two or more explanation fragments. Since the tool did not provide support for merging fragments, it did so by copying the explanation from one fragment to the end of the other fragment and removing the obsolete fragment. In only 10 cases, instructors attempted to create new explanations from scratch, but in the end

	A1	A2	A3	A4	A5	Total
Examples Created	6	2	2	1	1	12
Generated Explanations	126	32	44	8	27	237
Lines of Code being Explained by ChatGPT	55	12	21	2	9	99
Explanations Excluded	18	2	4	0	0	24
Explanations Liked	6	0	9	0	0	15
Explanations Edited	29	0	11	0	26	66
Explanations Removed	8	0	15	0	0	23

Table 4

Analysis of ChatGPT used explanations: Total count of generated, excluded, liked, and explained lines of code across 13 instances where the instructor added explanations to examples.

these explanations were removed. In other words, all remaining explanation fragments were originally generated by ChatGPT with some of them being edited later by the instructors. Apparently, the instructors preferred to edit the explanation fragments rather than create them from scratch. In total, the instructors edited 66 (27.84%) of ChatGPT generated explanation fragments, on average 1.4 times (stdev=0.55). Feedback from instructors indicated that most of their edits involved summarizing, adding missing details, or removing unnecessary parts. Table 4 shows that almost half of the generated fragments were used without being touched, saving a noticeable amount of instructor time.

	A1	A2	A3	A4	A5	Total
ChatGPT Edited Explanations	29	0	11	0	26	66
ChatGPT Explanation Edits	42	0	12	0	39	93
Average Levenshtein Ratio across all Final and Original ChatGPT Explanations	0.435	1	0.833	1	0.412	Average 0.736

Table 5

ChatGPT-generated explanations edits: Number of edits made by instructors to ChatGPT-generated explanations, along with a measure of similarity between the original and final edited version.

The average Levenshtein edit ratio for ChatGPT-generated explanations (edited and unedited) is 0.73 (Table 5), indicating a high acceptance rate for generated explanations. This indicator, however, is somewhat misleading since a portion of ChatGPT-generated explanations were edited because the first version of the tool evaluated in the study didn't provide direct support for reordering and merging the explanations, resulting in copy-pasting the explanations (as reported by A5 for whom the ratio dropped to 0.412). The table also points out that WEAT was able to support different editing approaches pursued by instructors. Some instructors spent more time reviewing the generated explanations before adding them to the example (A1), some prefer adding them to the examples and then evaluating and editing them (A3), while some used the generated explanations without changes.

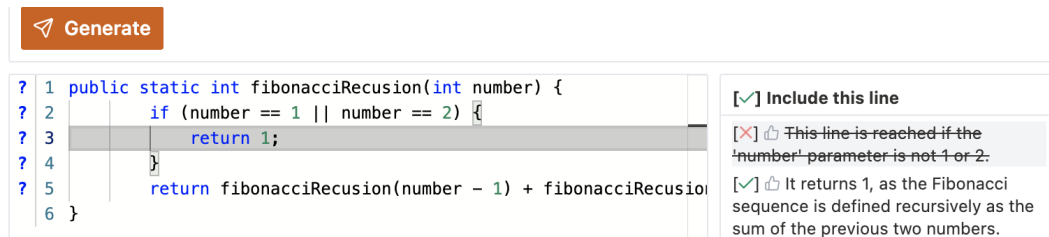


Figure 4: An incorrect explanation fragment generated by ChatGPT and excluded by the author (line 3).

6. Conclusion

In this paper, we introduce a worked code example authoring tool WEAT that supports human-AI co-creation in the process of developing such examples. WEAT supports human authors by using ChatGPT for the generation of line-by-line code explanations and by providing an interface to integrate this functionality into a balanced authoring process. To the best of our knowledge, this is the first attempt to develop an authoring tool that produces worked examples through human-AI collaboration.

To develop WEAT, we performed several rounds of feasibility studies. These studies supported the need for a human-AI co-creation in authoring worked examples. As the studies showed, in the majority of cases, the explanations generated by ChatGPT with a carefully tuned prompt were positively evaluated by authors and students. However, in a good fraction of cases they were inferior to the explanations provided by experts. The study also revealed that on average experts can create explanations that are more easily readable and closer to the explanations generated by the students themselves. With this data, we hypothesized that human-AI co-creation could offer the “best of both worlds” solution where good explanations could be simply accepted by authors, while inferior or hard-to-understand explanations could be improved.

An evaluation of WEAT system with five course instructors supported these expectations and provided strong evidence in favor of co-creation. As the log analysis demonstrated, in many cases, instructors choose to accept generated explanations without changes, which should have decreased the time and effort required for example creation. Yet in other cases, the instructor rejected or edited the generated explanation to achieve the desired quality. In some cases, explanations were rejected by being simply incorrect, which stresses the importance of human presence in the authoring process. The interview with authors revealed several cases where authors acted inefficiently due to specific interface issues, such as blocking the main edit window by the generation dialog or the lack of tools to move or merge fragments. Now we are using these observations to develop an improved version of WEAT.

As the first step towards this important goal, our work has limitations. Most importantly, the scale of our evaluation is relatively small. Since we targeted real instructors as users in our evaluation process, we were able to recruit only five qualified subjects. Additionally, since the study was done at the beginning of the semester when instructors were busy setting up their classes, they created only 12 examples using this tool. To obtain more reliable data, we plan a larger-scale semester-long study by engaging instructors to create a variety of worked examples of varying difficulty and use them in their classes. Such a study will also enable us to

assess the quality of explanations produced through human-AI collaboration and their value for students in introductory programming classes.

References

- [1] M. C. Linn, M. J. Clancy, The case for case studies of programming problems, *Commun. ACM* 35 (1992) 121–132.
- [2] H. M. Deitel, P. J. Deitel, *C How to Program*, 2nd Edition, Prentice Hall, New York, 1994.
- [3] A. Kelley, I. Pohl, *C by Dissection : The Essentials of C Programming*, Addison-Wesley, New York, 1995.
- [4] P. Brusilovsky, M. V. Yudelson, I.-H. Hsiao, Problem solving examples as first class objects in educational digital libraries: Three obstacles to overcome, *Journal of Educational Multimedia and Hypermedia* 18 (2009) 267–288.
- [5] R. Sharrock, E. Hamonic, M. Hiron, S. Carlier, Codecast: An innovative technology to facilitate teaching and learning computer programming in a c language online course, *Proceedings of the Fourth (2017) ACM Conference on Learning @ Scale (2017)*.
- [6] K. Khandwala, P. J. Guo, Codemotion: expanding the design space of learner interactions with computer programming tutorial videos, *Proceedings of the Fifth Annual ACM Conference on Learning at Scale (2018)*.
- [7] J. Park, Y. H. Park, J. Kim, J. Cha, S. Kim, A. H. Oh, Elicast: embedding interactive exercises in instructional programming screencasts, *Proceedings of the Fifth Annual ACM Conference on Learning at Scale (2018)*.
- [8] R. Hosseini, K. Akhuseyinoglu, P. Brusilovsky, L. Malmi, K. Pollari-Malmi, C. Schunn, T. Sirkiä, Improving engagement in program construction examples for learning python programming, *International Journal of Artificial Intelligence in Education* 30 (2020) 299–336.
- [9] I.-H. Hsiao, P. Brusilovsky, The role of community feedback in the student example authoring process: an evaluation of annotex, *British Journal of Educational Technology* 42 (2011) 482–499.
- [10] A. Davidovic, J. R. Warren, E. Trichina, Learning benefits of structural example-based adaptive tutoring systems, *IEEE Trans. Educ.* 46 (2003) 241–251.
- [11] B. B. Morrison, L. E. Margulieux, B. Ericson, M. Guzdial, Subgoals help students solve parsons problems, *Proceedings of the 47th ACM Technical Symposium on Computing Science Education (2016)*.
- [12] B. Ericson, M. Guzdial, B. B. Morrison, Analysis of interactive features designed to enhance learning in an ebook, *Proceedings of the eleventh annual International Conference on International Computing Education Research (2015)*.
- [13] M. T. H. Chi, J. Adams, E. B. Bogusch, C. Bruchok, S. Kang, M. Lancaster, R. Levy, N. Li, K. L. McEldoon, G. S. Stump, R. Wylie, D. Xu, D. L. Yaghmourian, Translating the icap theory of cognitive engagement into practice, *Cognitive Science* 42 (2018) 1777–1832.
- [14] J. Zamfirescu-Pereira, R. Y. Wong, B. Hartmann, Q. Yang, Why johnny can't prompt: How non-ai experts try (and fail) to design llm prompts, in: *Proceedings of the 2023 CHI*

- Conference on Human Factors in Computing Systems, CHI '23, Association for Computing Machinery, New York, NY, USA, 2023.
- [15] S. MacNeil, A. Tran, A. Hellas, J. Kim, S. Sarsa, P. Denny, S. Bernstein, J. Leinonen, 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, SIGCSE 2023, Association for Computing Machinery, New York, NY, USA, 2023, p. 931–937.
 - [16] J. Leinonen, P. Denny, S. MacNeil, S. Sarsa, S. Bernstein, J. Kim, A. Tran, A. Hellas, Comparing code explanations created by students and large language models, 2023.
 - [17] J. Li, S. Tworkowski, Y. Wu, R. Mooney, Explaining competitive-level programming solutions using llms, 2023.
 - [18] E. Chen, R. Huang, H.-S. Chen, Y.-H. Tseng, L.-Y. Li, Gptutor: A chatgpt-powered programming tool for code explanation, in: N. Wang, G. Rebolledo-Mendez, V. Dimitrova, N. Matsuda, O. C. Santos (Eds.), Artificial Intelligence in Education. Posters and Late Breaking Results, Workshops and Tutorials, Industry and Innovation Tracks, Practitioners, Doctoral Consortium and Blue Sky, Springer Nature Switzerland, Cham, 2023, pp. 321–327.
 - [19] S. Sarsa, P. Denny, A. Hellas, J. Leinonen, 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, ICER '22, Association for Computing Machinery, New York, NY, USA, 2022, p. 27–43.
 - [20] H. Tian, W. Lu, T. O. Li, X. Tang, S.-C. Cheung, J. Klein, T. F. Bissyandé, Is chatgpt the ultimate programming assistant – how far is it?, 2023.
 - [21] D. Zhou, N. Scharli, L. Hou, J. Wei, N. Scales, X. Wang, D. Schuurmans, O. Bousquet, Q. Le, E. H. Hsin Chi, Least-to-most prompting enables complex reasoning in large language models, ArXiv (2022).
 - [22] J. White, Q. Fu, S. Hays, M. Sandborn, C. Olea, H. Gilbert, A. Elnashar, J. Spencer-Smith, D. C. Schmidt, A prompt pattern catalog to enhance prompt engineering with chatgpt, 2023.
 - [23] M. Hassany, P. Brusilovsky, J. Ke, K. Akhuseyinoglu, A. B. Lekshmi Narayanan, Authoring Worked Examples for Java Programming with Human-AI Collaboration, Report arXiv:2312.02105, arXiv, 2023. URL: <https://doi.org/10.48550/arXiv.2312.02105>.
 - [24] A.-B. Lekshmi-Narayanan, P. Oli, J. Chapagain, M. Hassany, R. Banjade, P. Brusilovsky, V. Rus, Explaining code examples in introductory programming courses: Llm vs humans, in: Workshop on AI for Education - Bridging Innovation and Responsibility at AAAI 2024, 2024.