

# Explaining Code Examples in Introductory Programming Courses: LLM vs Humans

Arun-Balajiee Lekshmi-Narayanan \*

Priti Oli \*

Jeevan Chapagain

Mohammad Hassany

Rabin Banjade

Peter Brusilovsky

Vasile Rus

*University of Pittsburgh, PA, USA. 15260*

*University of Memphis, TN, USA, 38152*

ARL122@PITT.EDU

POLI@MEMPHIS.EDU

JCHPGAIN@MEMPHIS.EDU

MOH70@PITT.EDU

RBNJADE1@MEMPHIS.EDU

PETERB@PITT.EDU

VRUS@MEMPHIS.EDU

## Abstract

Worked examples, which present an explained code for solving typical programming problems 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 explanations for many examples typically used in a programming class. In this paper, we assess the feasibility of using LLMs to generate code explanations for passive and active example exploration systems. To achieve this goal, we compare the code explanations generated by chatGPT with the explanations generated by both experts and students.

**Keywords:** Programming, Worked Examples, Code Explanations, ChatGPT

## 1. Introduction

Program code examples (known also as worked examples) play a crucial role in learning how to program (Linn and Clancy, 1992). Instructors use examples extensively to demonstrate the semantics of the programming language being taught and to highlight the fundamental coding patterns. Programming textbooks allocate considerable space to present and explain code examples. To make the process of studying code examples more interactive, CS education researchers developed a range of tools to engage students in the study of code examples. These tools include codecasts (Sharrock et al., 2017), interactive example explorers (Hosseini et al., 2020), and tutoring systems (Oli et al., 2023).

An important component in all types of program examples is code explanations associated with code lines or chunks. The explanations connect examples with general programming knowledge explaining the role and function of code fragments or their behavior. In textbooks, these explanations are usually presented as comments in the code or as explanations on the margins. The example explorer tools allow students to examine these explanations interactively (Hosseini et al., 2020). Tutoring systems, which engage students in explaining the code, use instructor explanations to assess student responses (Chapagain et al., 2022) and provide scaffolding (Oli et al., 2023). The explanations must be *authored*

---

\* Both authors contributed equally.

by instructors or domain experts, i.e., written and integrated into a specific system. As the experience of the last 10 years demonstrated, these explanations are hard to obtain. Being enthusiastic about sharing *the code* of their worked examples with others, instructors generally do not have time or patience to properly author *explanations* of their examples. Indeed, creating just one explained example could take 30 minutes even in the presence of authoring tools (Hosseini et al., 2020; Sharrock et al., 2017). As a result, the volume of worked examples available to students in a typical introductory programming class is low.

To address this *authoring bottleneck*, researchers explored *learner-sourcing*, that is, engaging students in the creation and review of explanations of instructor-provided code (Hsiao and Brusilovsky, 2011) and the automatic extraction of explanations from lecture recordings (Khandwala and Guo, 2018). In this paper, we explore the feasibility of human-AI collaboration in creating explained code examples. With this approach, the instructor provides the code for their favorite examples. The AI engine based on large language models (LLM) examines the example code and generates explanations for each code line. The explanations are reviewed and, if necessary, edited by the instructor.

To assess the feasibility of this approach, it is important to compare the code explanations produced by LLMs such as ChatGPT with explanations produced by humans. To use ChatGPT explanations in example explorer systems, we need to check how similar they are in language, semantics, and style used to code explanations produced by instructors and domain experts. To use ChatGPT explanations to assess student responses in tutoring systems, we need to assess how close these explanations are to the explanations produced by students when explaining the code. In this paper, we employ a range of analytical approaches to compare ChatGPT code explanations with explanations produced by both experts and students. Following a review of past work on using LLM for code explanations, we present the method and datasets used in the study and review the results.

## 2. Related Work: Use of LLMs for Code Explanations

Multiple researchers have explored code summarization (Phillips et al., 2022) and explanations using transformer models (Choi et al., 2023; Peng et al., 2022), abstract syntax trees (Shi et al., 2022), and Tree-LSTM (Tian et al., 2023). With the announcement of ChatGPT, several research teams explored the use of LLM for code explanations using ChatGPT 3 (Zamfirescu-Pereira et al., 2023; MacNeil et al., 2023; Leinonen et al., 2023), GPT 3.5 (MacNeil et al., 2023; Li et al., 2023; Chen et al., 2023), GPT 4 (Li et al., 2023), OpenAI Codex (Sarsa et al., 2022; Tian et al., 2023; MacNeil et al., 2023), and GitHub Copilot (Chen et al., 2023). These LLMs were used to generate explanations at different levels of abstraction (line-by-line, step-by-step, and high-level summary). Sarsa et al. (2022) observed that ChatGPT can generate better explanations at line-by-line level.

The explanations and summaries generated by these LLMs were mostly evaluated by authors (Sarsa et al., 2022), students (MacNeil et al., 2023; Leinonen et al., 2023), and tool users (Chen et al., 2023). Most recently, attempts have been made to compare code explanations generated by humans and LLMs. Sarsa et al. (2022) reported that students rated LLM-generated explanations as useful, easier, and more accurate than learner-sourced explanations (Leinonen et al., 2023). In this work, we attempted a more formal approach to compare the line-by-line code explanations generated by ChatGPT, students, and experts

using a range of quantitative metrics. Our goal was to generate insights into the use of ChatGPT code explanations in the learning process.

### 3. Dataset Collection

To produce the evaluation dataset, we collected line-by-line explanations from three types of sources - experts, ChatGPT, and students - for four Java worked examples selected from different topics of the PCEX example exploration system (Hosseini et al., 2020). The examples represent programming problems of different difficulty levels. The simpler program involved array search and print statements, while the hardest program focused more on object-oriented principles. The selected examples were used in multiple Java classes and included line-by-line explanations produced by the instructors. The summary of the types of explanation collected is shown in Table 1 and the collection process is explained below. Figure 1 shows sample explanations from different sources.

*Expert Explanations:* We used one set of expert explanations available at PCEX and collected the second set of line-by-line explanations from different experts.

<p><b>Program:</b> PointTester; <b>Line number:</b> 12; <b>Line code:</b> <code>private int y;</code></p> <p><b>Expert1:</b> Every object of the Point class will have its own y-coordinate. Therefore, we need to declare an instance variable for the class to store the y-coordinate of the point. We declare it as int because we want to have integer coordinates for the point. Note that an instance variable is a variable defined in a class, for which each instantiated object of the class has a separate copy, or instance.</p> <p><b>Expert2:</b> The instance variables are declared as private to prevent direct access to them from outside the class. In this way, no unexpected modifications to a Point object's data are possible.</p> <p><b>S:</b> This line declares a private integer variable named "y" to store the y coordinate of a point.</p> <p><b>A:</b> This line defines a private instance variable 'y' of type int in the Point class. It contributes directly to the program's objective of storing the y-coordinate of the point.</p> <p><b>E:</b> This line declares a private instance variable 'y' of type int in the Point class to store the y-coordinate of a point. Declaring the y-coordinate variable is essential for keeping track of the point's position and contributes directly to the program's objective of storing the point's coordinates.</p> <p><b>Student1:</b> initialize a private value inside the point class with no value yet</p> <p><b>Student2:</b> Declares the private int variable y.</p> <p><b>Student3:</b> Creates a private int that can only be accessed by class Point called int y</p> <p>...</p> <p><b>Student59:</b> private variable used to store the value entered into the value of the y coordinate.</p>
--

Figure 1: Sample Data illustrating the various entries under each column in our dataset.

*Student Explanations:* We performed a user study in which students of a Java programming course were asked to write explanations for each line of the code examples selected for

Explanation Type	N	Definition
Experts	2	Source Code Line-by-Line Explanations by Experts
Students	60	Source Code Line-by-Line Explanations by Students
Simple Prompt (S)	4	ChatGPT Explanations with simple prompt (Section 3)
Advanced Prompt (A)	4	ChatGPT Explanations with advanced prompt (Figure 2)
Extended Prompt (E)	4	ChatGPT Explanations with extended prompt (Figure 2)
ChatGPT	12	Aggregated representation for all ChatGPT prompts

Table 1: A summary of explanation sources used in the study.

the study. In total, we collected line-by-line explanations from 60 students. For example, for the line of code “*private int y*” a student participating in the study explained “*Creates a new object class called Point*”.

**ChatGPT Explanations:** We performed a sequence of internal studies to build ChatGPT prompts, which can produce the most useful line-by-line example explanations. For the final evaluation, we selected three prompts that produced good explanations on three different levels of detail. We used gpt-3.5-turbo to generate four sets of line-by-line explanations for each selected example using each of these prompts. To increase the diversity of explanations, the first set was generated with temperature value 0 and three additional sets were generated with temperature value 1 each time clearing the history.

**Simple Prompt:** The *Simple prompt* used in the study included the code of the worked example and the following instruction: “*Provide a line-by-line self-explanation for each line of code in the Java program above*”. The explanations generated by this prompt are concise and elicit the goal of each line of code quite well.

**Advanced Prompt:** The *Advanced prompt* used both the problem statement and the example code along with more elaborate instructions for ChatGPT. The role of “a professor who teaches computer programming” is assigned to the system to provide a context. The prompt also asked ChatGPT to provide reasons why the line needs to be explained. We observed that ChatGPT cannot always associate line numbers correctly, so each line in the program source code was annotated with its associated line number. An output format was defined to process the results digestible by our automation script. Prompt details are provided in Figure 2 (Iteration #1 ).

**Extended Prompt:** To obtain the most elaborate ChatGPT explanations, we used *Extended prompt*, which requested ChatGPT to further enhance the explanations generated by the Advanced prompt (Iteration #2 in Figure 2), with a focus on consistency and coverage of the generated content.

#### 4. Evaluation Metrics

**Lexical Metrics:** We report the *lexical diversity* and *lexical density* of the generated explanations to assess the richness, informativeness, and conciseness of the generated text (Johansson, 2008). *Lexical diversity* is the range of variety of distinct words or vocabulary used within a specific text. *Lexical density* refers to the measure of the variety of different lexical words present in a text, including nouns, adjectives, verbs, and adverbs, which collectively contribute to the overall meaning of the text. A recent paper considers the use of

**Role (System):**  
 You are a professor who teaches computer programming.

---

**Role (User) - Iteration #1:**  
 Given the following program description and accompanying source code, identify and explain lines of the code that contributes directly to the program objectives and goals.  
  
 When considering each identified line, ensure explanations provide the reasons that led to the line inclusion, prioritizing them based on their relative importance while also preventing any unnecessary duplication or repetition of information.  
  
*Program Description:*  
**[program description]**  
  
*Program Source Code:*  
 The line number is defined as /\*line\_num\*/ at the start of each line.  
  

```

/*java
/*1*/[program lines]
*/

```

  
*Output format:*  
 Reply ONLY with a JSON array where each element, representing a "line of code," includes "line\_num" and an "explanations" array. For example:  

```

/*json
[ { "line_num": "2", "explanations": [ "explanation ...", "explanation ...", ... ] }, ... ]
*/

```

---

**Role (User) - Iteration #2:**  
 Update your explanations with more insightful and complementary YET COMPLETELY new explanations. If you missed a line, this is the time to include them.

Figure 2: This ChatGPT Prompt template considers the case that ChatGPT could generate a better explanation with an additional “nudge” as observed above. In most cases, the generated explanations using the prompt at the second iteration produces richer explanation than the first iteration.

lexical diversity as a metric (Cegin et al., 2023) to compare content generated by humans and ChatGPT. We also report on the total number of tokens on each explanation to provide more insight into the comparison of the lexical features across each source.

*Readability Metrics:* We consider 3 popular metrics (Denny et al., 2020, 2021), namely, Flesch-Kincaid Grade Level, Gunning Fog and Flesch Reading Ease. These metrics esti-

mate the grade level expectation of the text such as the years of formal education required to read the text and the ease of reading a piece of text respectively. We use the TextDescriptives (Hansen et al., 2023) package in Python<sup>1</sup> to calculate these scores.

*Similarity Metrics:* We use four evaluation metrics: Character-based metric chrF (Popović, 2015), Word-based metric METEOR (Banerjee and Lavie, 2005), and Embedding-based metrics BERTScore (Zhang et al., 2019) and Universal Sentence Encoder (USE) (Cer et al., 2018) to compare explanation generated by different sources. chrF (character n-gram F-score) measures the character-level matching between the reference text and the machine-generated text considering both precision and recall. METEOR considers the similarity between words and assesses word overlap between the two texts. BERTScore is an automated evaluation metric for text generation that assesses the similarity between candidate and reference sentences by comparing the contextual embeddings of individual tokens using cosine similarity. USE is a transformer-based model that transforms text into high-dimensional vectors, enabling the computation of similarity between two texts based on their vector representations. Haque et al. (2022) and Roy et al. (2021) have pointed out that METEOR, chrF (Popović, 2015), and USE (Cer et al., 2018) metrics are better aligned with human preferences of code summarization as these metrics assign partial credits to words. We also use BertScore to evaluate the generated explanation primarily due to its extensive use as a reliable measure for evaluating the faithfulness of LLMs (Ji et al., 2023). Consequently, traditional metrics like BLEU (Papineni et al., 2002) which solely rely on word overlap, are now considered outdated and are not included in our reporting.

## 5. Results

We compared the explanations using metrics aggregated over all 33 explainable lines of four examples generated by each expert, each student, and each round of ChatGPT generation. Table 2 reports the medians for each source *type*.

*Lexical Metrics:* As Table 2 shows, explanations produced by experts and ChatGPT are more than twice as long as explanations produced by students (as measured by the number of tokens). Students also use considerably fewer unique words in their explanations (lexical diversity) hinting that their vocabulary is more narrow than the vocabulary of experts and ChatGPT. The length and lexical diversity of explanations generated by ChatGPT and experts varied, with *Simple prompt* generating the shortest, *Extended prompt* the longest explanations, and expert explanations positioned between *Advanced* and *Extended* prompts. An ANOVA analysis of the lexical diversity of the explanation generated by experts, ChatGPT, and students indicated statistically significant variations among the groups (F-statistic = 25.07,  $p < 0.05$ ). The data also show that the explanations produced by the students are not only shorter than those by experts and ChatGPT, but also have a higher lexical density, suggesting that the students explain the code in a more “concentrated” way. However, the ANOVA analysis of the lexical density does not indicate any significant difference (F-statistic = 2.5,  $p = 0.08$ ) between the explanations in terms of lexical density.

*Readability Metrics:* One-way ANOVA revealed that the Gunning-Fog readability scores are significantly different between explanation sources (experts, students, and ChatGPT)

1. <https://hlasse.github.io/TextDescriptives/readability.html>



Source	N	Vocabulary	Lexical Density	# of Tokens	GF	FRE	FK
Experts	2	209.0	0.48	690.0	8.46	78.45	6.18
S	4	165.0	0.45	517.5	8.67	82.34	6.35
A	4	185.5	0.48	625.0	9.91	72.63	7.15
E	4	238.0	0.49	769.5	11.09	69.64	7.83
ChatGPT	12	179.5	0.48	625.0	8.99	75.41	6.69
Students	60	116.5	0.54	249.5	8.02	80.48	5.62

Table 2: Median lexical and readability metrics for different sources of explanations (FRE = Flesch-Reading Ease, FK = Flesch-Kincaid, GF = Gunning Fog)

( $p < 0.001$ ). From post hoc comparisons, the most significant differences were observed between explanations produced by *Extended prompt* and experts ( $p = 0.0102$ ) as well as students ( $p < 0.0001$ ). No significance was observed between experts and students ( $p = 0.0848$ ). Overall, for all metrics, the ChatGPT explanations are relatively less readable than those of experts, which are less readable than those of students (Table 2).

*Similarity Metrics:* We applied similarity metrics to calculate the similarity between the explanations provided by ChatGPT, students and experts for each line of code. We computed pairwise similarity scores between ChatGPT, Expert, and Student explanations averaged over all the 33 lines from all 4 programs. Table 3 shows that the explanations generated by ChatGPT exhibit a consistently higher average similarity score to the expert explanations compared to those generated by students. This suggests that the explanations generated by ChatGPT are more closely aligned with expert explanations than with student explanations across all metrics.

Mann-Whitney U-tests indicate significantly higher alignment between ChatGPT and expert explanations than between student and expert explanations (F-statistic = 48.0,  $p < 0.05$  for METEOR, F-statistic = 205.0,  $p < 0.05$  for USE and F-statistics=288.0,  $p < 0.05$  for BERTScore). We also observed that *Simple* prompt (S) generated explanations aligned more closely with expert explanations than *Advanced* and *Extended* prompts (A and E). The higher semantic alignment between the expert explanation and the ChatGPT explanations generated with a *Simple* rather than *Advanced* prompt could be explained by the nature of the *Advanced* prompt, which was specifically engineered to produce very detailed explanations which experts rarely have time to produce. In particular, explanations generated using *Advanced* prompt frequently explained why each line of the code is important, while neither expert explanations nor explanations generated with *Simple* prompt consistently explained the importance of each line in the code.

## 6. Conclusions and Future Work

In this work, our goal was to assess the feasibility of using ChatGPT to generate line-by-line code explanations to be used in worked-out examples in place of currently used expert explanations. To achieve this goal, we compared the ChatGPT explanations generated by different prompts with the student and expert explanations using a range of metrics. Our results indicate that the explanations generated by ChatGPT are lexically and semantically

Reference	Source	chrF	METEOR	USE	BERTScore
Expert	ChatGPT(S)	0.32	0.28	0.54	0.89
Expert	ChatGPT(A)	0.31	0.27	0.49	0.709
Expert	ChatGPT(E)	0.32	0.28	0.48	0.712
Expert	ChatGPT (All)	0.32	0.27	0.50	0.75
Expert	Student	0.33	0.144	0.33	0.63
ChatGPT(S)	Student	0.22	0.21	0.34	0.60
ChatGPT(A)	Student	0.18	0.150	0.254	0.450
ChatGPT(E)	Student	0.18	0.151	0.255	0.458
ChatGPT(All)	Student	0.19	0.17	0.28	0.50
ChatGPT(S)	ChatGPT(A)	0.33	0.30	0.50	0.72
ChatGPT(S)	ChatGPT(E)	0.32	0.28	0.50	0.73
ChatGPT(A)	ChatGPT(E)	0.44	0.43	0.56	0.69

Table 3: Assessing alignment (larger is better) between sources of explanations

similar to the explanations generated by experts and could potentially resolve the authoring bottleneck. However, their lower readability level might be an obstacle for less-prepared students. We also observed a considerable difference between the explanations produced by students and explanations produced by experts and ChatGPT, which might affect the efficiency of both sources of explanation in active example tutors where student explanations are assessed by comparing them with expert explanations. Both issues require a deeper investigation, which we plan to perform through a user study.

## 7. Acknowledgement

This work has been supported by NSF awards 1822816 and 1822752. The opinions, findings, and results are solely those of the authors and do not reflect those of NSF.

## References

- Satanjeev Banerjee and Alon Lavie. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, 2005.
- Ján Cegin, Jakub Simko, and Peter Brusilovsky. Chatgpt to replace crowdsourcing of paraphrases for intent classification: Higher diversity and comparable model robustness. *ArXiv*, abs/2305.12947, 2023.
- Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St John, Noah Constant, Mario Guajardo-Cespedes, Steve Yuan, Chris Tar, et al. Universal sentence encoder. *arXiv preprint arXiv:1803.11175*, 2018.



- Jeevan Chapagain, Lasang Tamang, Rabin Banjade, Priti Oli, and Vasile Rus. Automated assessment of student self-explanation during source code comprehension. In *The International FLAIRS Conference Proceedings*, 2022.
- Eason Chen, Ray Huang, Han-Shin Chen, Yuen-Hsien Tseng, and Liang-Yi Li. GPTutor: A ChatGPT-powered programming tool for code explanation. In *Artificial Intelligence in Education.*, 2023.
- Yunseok Choi, Cheolwon Na, Hyojun Kim, and Jee-Hyong Lee. Readsum: Retrieval-augmented adaptive transformer for source code summarization. *IEEE Access*, 11:51155–51165, 2023.
- Paul Denny, James Prather, and Brett A. Becker. Error message readability and novice debugging performance. In *Proceedings of the 2020 ACM Conference on Innovation and Technology in Computer Science Education*, 2020.
- Paul Denny, James Prather, Brett A. Becker, Catherine Mooney, John Homer, Zachary C Albrecht, and Garrett B. Powell. On designing programming error messages for novices: Readability and its constituent factors. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, 2021.
- Lasse Hansen, Ludvig Renbo Olsen, and Kenneth C. Enevoldsen. Textdescriptives: A python package for calculating a large variety of metrics from text. *J. Open Source Softw.*, 8:5153, 2023.
- Sakib Haque, Zachary Eberhart, Aakash Bansal, and Collin McMillan. Semantic similarity metrics for evaluating source code summarization. In *Proceedings of the 30th IEEE/ACM International Conference on Program Comprehension*, 2022.
- Roya Hosseini, Kamil Akhuseyinoglu, Peter Brusilovsky, Lauri Malmi, Kerttu Pollari-Malmi, Christian Schunn, and Teemu Sirkiä. Improving engagement in program construction examples for learning python programming. *International Journal of Artificial Intelligence in Education*, 2020.
- I-Han Hsiao and Peter Brusilovsky. The role of community feedback in the student example authoring process: an evaluation of annotex. *British Journal of Educational Technology*, 2011.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 2023.
- Victoria Johansson. Lexical diversity and lexical density in speech and writing: A developmental perspective. *Working papers/Lund University, Department of Linguistics and Phonetics*, 53:61–79, 2008.
- Kandarp Khandwala and Philip 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.

- Juho Leinonen, Paul Denny, Stephen MacNeil, Sami Sarsa, Seth Bernstein, Joanne Kim, Andrew Tran, and Arto Hellas. *Comparing Code Explanations Created by Students and Large Language Models*. 2023. arXiv:2304.03938.
- Jierui Li, Szymon Tworkowski, Yingying Wu, and Raymond Mooney. Explaining competitive-level programming solutions using LLMs, 2023. arXiv:2307.05337.
- Marcia C. Linn and Michael J. Clancy. The case for case studies of programming problems. *Commun. ACM*, 1992.
- Stephen MacNeil, Andrew Tran, Arto Hellas, Joanne Kim, Sami Sarsa, Paul Denny, Seth Bernstein, and Juho 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*, 2023.
- Priti Oli, Rabin Banjade, Arun Balajjee Lekshmi Narayanan, Jeevan Chapagain, Lasang Jimba Tamang, Peter Brusilovsky, and Vasile Rus. Improving code comprehension through scaffolded self-explanations. In *Proceedings of 24th International Conference on Artificial Intelligence in Education*, 2023.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, 2002.
- Han Peng, Ge Li, Yunfei Zhao, and Zhi Jin. Rethinking positional encoding in tree transformer for code representation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, 2022.
- Jesse Phillips, David Bowes, Mahmoud El-Haj, and Tracy Hall. Improved evaluation of automatic source code summarisation. *Proceedings of the 2nd Workshop on Natural Language Generation, Evaluation, and Metrics (GEM)*, 2022.
- Maja Popović. chrF: character n-gram f-score for automatic mt evaluation. In *Proceedings of the tenth workshop on statistical machine translation*, 2015.
- Devjeet Roy, Sarah Fakhoury, and Venera Arnaoudova. Reassessing automatic evaluation metrics for code summarization tasks. In *Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 2021.
- 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 2022 ACM Conference on International Computing Education Research*, 2022.
- Rémi Sharrock, Ella Hamonic, Mathias Hiron, and Sébastien Carlier. Codecast: An innovative technology to facilitate teaching and learning computer programming in a c language online course. *Proceedings of the Fourth ACM Conference on Learning @ Scale*, 2017.
- Yang Shi, Min Chi, Tiffany Barnes, and Thomas William Price. Code-DKT: A code-based knowledge tracing model for programming tasks, 2022. ArXiv:2206.03545.

Haoye Tian, Weiqi Lu, Tsz On Li, Xunzhu Tang, Shing-Chi Cheung, Jacques Klein, and Tegawendé F. Bissyandé. Is ChatGPT the ultimate programming assistant – how far is it?, 2023. arXiv:2304.11938.

J.D. Zamfirescu-Pereira, Richmond Y. Wong, Bjoern Hartmann, and Qian 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*, 2023.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*, 2019.