



On the Opportunities of Large Language Models for Programming Process Data

John Edwards
Computer Science
Utah State University
Utah, USA
john.edwards@usu.edu

Arto Hellas
Aalto University
Espoo, Finland
arto.hellas@aalto.fi

Juho Leinonen
Aalto University
Espoo, Finland
juho.2.leinonen@aalto.fi

Abstract

Computing educators and researchers have long used programming process data to understand how students construct programs and address challenges. Despite its potential, fully automated feedback systems remain underexplored. The emergence of Large Language Models (LLMs) offers new opportunities for analyzing programming data and providing formative feedback. This study explores using LLMs to summarize programming processes and deliver formative feedback. A case study analyzed keystroke-level data from an introductory programming course, processed into code snapshots. Three state-of-the-art LLMs – Claude 3 Opus, GPT-4 Turbo, and LLaMa2 70B Chat – were evaluated for their feedback capabilities. Results show LLMs effectively provide tailored feedback, emphasizing incremental development, algorithmic planning, and code readability. Our findings highlight the potential of combining keystroke data with LLMs to automate formative feedback, showing that the computing education research and practice community is again one step closer to automating formative programming process feedback.

CCS Concepts

• **Social and professional topics** → **Computing education.**

Keywords

programming process data, large language models, generative AI, programming process feedback, programming process summarization, keystroke data

ACM Reference Format:

John Edwards, Arto Hellas, and Juho Leinonen. 2025. On the Opportunities of Large Language Models for Programming Process Data. In *27th Australasian Computing Education Conference (ACE 2025), February 12–13, 2025, Brisbane, QLD, Australia*. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3716640.3716652>

1 Introduction

Feedback can have a tremendous impact on learning and achievement [28]. The level of detail of the feedback influences its effectiveness [87], and feedback can be given at many levels ranging from

targeting how to work on and complete specific tasks to considering personal characteristics and behavior [28, 38, 65]. In teaching and learning programming, automated assessment systems have been a key tool for providing feedback at a scale already for more than a half a century [32, 38, 67]. Researchers have sought to automate step-by-step guidance [85], provide hints during the programming process [61], improve programming error messages [6], and aid in providing textual feedback by grouping similar code submissions together [25, 39, 64].

To support the understanding of how novices construct programs, researchers and educators have been collecting increasing amounts of data from students' programming process [33]. Such data can be collected at multiple granularities, ranging from final course assignment submissions to individual keystrokes from solving the assignments [33]. Programming process data has been, for example, used to play back how students construct their programs step by step or keystroke by keystroke to create a broader understanding of the process [29, 44, 80, 90]. So far, despite shared efforts towards providing timely feedback to students [35], the potential of fine-grained programming process data for feedback purposes is still largely untapped.

Large Language Models (LLMs) are a potential tool for realizing the transformation of programming process data into actionable feedback items. Within Computing Education Research, LLMs have broadened the horizon of what computing education researchers and practitioners can achieve [71], calling even for rethinking how computer science and programming is taught [17, 18]. Large Language Models have been shown to help in creating assignments [56, 79], improve error messages [51, 78], fix students' code [42], explain code [50, 59, 60], and respond to help requests [11, 31, 40]. At their core, LLMs are tools that allow transforming text-based content from one form to another, drawing on the data that they have been trained with in the process and the instructions provided by the user [66, 83].

In this article, we discuss the potentials of LLMs for programming process data. We outline a case study of using LLMs for analyzing programming process data. In the case study, we first use LLMs for summarizing the programming process, followed by asking LLMs to provide feedback on the process. This article is structured as follows. In Section 2, we outline prior research on programming process data. Building on Section 2, Section 3 outlines our vision for the possibilities of LLMs for programming process data. In Section 4, we outline our case study where we leveraged LLMs for analyzing programming process data. Finally, in Section 5, we summarize our discussion and outline possible future directions for research on using LLMs for programming process analytics.



This work is licensed under a Creative Commons Attribution International 4.0 License.

ACE 2025, Brisbane, QLD, Australia

© 2025 Copyright held by the owner/author(s).

ACM ISBN 979-8-4007-1425-2/25/02

<https://doi.org/10.1145/3716640.3716652>

2 Programming Process Data

In this section, we discuss the types and different uses of programming process data. Two surveys [20, 33] are resources for discussion of both the different granularities of raw process data and the wide variety of studies and research questions being pursued using this type of data.

2.1 Types of Data

Programming process data has been collected in various forms. Ihanola et al. [33] discuss process data at six discrete levels and Edwards et al. [20] define four levels, which are, starting from the least granular: submissions and commits [e.g., 2, 41]; executions, compilations, and file saves [e.g., 14, 34, 36]; line-level edits, which capture all contiguous changes in a single line in a single event [e.g., 9, 12]; and finally, keystroke and character edits [e.g., 13, 20, 48, 55]. Our interest is in using occasional snapshots of code and analyze the progression of one snapshot to another. Any granularity of process data would potentially be suitable for our purposes, but, in practice, students may not submit or commit their data often enough, nor do executions or compilations guarantee sufficient temporal resolution. However, line-level and keystroke-level data are generally too fine-grained. In this paper, we use keystroke-level data but discretize it into code snapshots that immediately precede a break in keystrokes of at least five minutes.

The five-minute threshold is based on a probabilistic model by Hart et al. [27] that shows that breaks of five minutes statistically indicate a 50% chance that the student has disengaged. Among the many other thresholds that have been used, which arbitrarily range from 60 seconds to one hour [37, 53, 63, 73], one study used a threshold of five minutes [45], but it pre-dated Hart’s statistical model and was arbitrarily chosen.

A small number of public datasets are available that would be suitable for our experiments, including Blackbox [9, 10, 12], which has line-level edit data, and a CSEDm challenge dataset¹ which has compile-level granularity. We chose the dataset provided by Edwards et al. [20], which is keystroke-level, because it gives us more flexibility over the temporal resolution of code snapshots.

2.2 Uses of Programming Process Data

Programming process data has been used in biometrics and authentication [43, 54, 55, 57, 62] and plagiarism detection [30, 77, 80], but the most common use of programming process data is in predicting course outcomes. Early prediction work used features of compiler errors, including frequency and repetition [1, 5, 15, 34, 86]. Later work incorporated other features derived from process data, including number of keystrokes [81], number of IDE hints [23, 24], time-on-task [47, 75], and many others [20].

Pertinent to our work presented in this paper are studies that looked at how the code itself is constructed, though most of these use rudimentary measures. They include measuring incremental test writing and execution to derive how often a student works on their program [36]; clustering of students into tinkerers and planners [8]; abstractly visualizing student progress toward a solution [70, 80]; and tools to interactively visualize code snapshots to facilitate discussion [89]. Intelligent Tutoring Systems (ITSs) have

used various features for machine learning techniques to determine hints and feedback as students write code [4, 72, 74, 76]. In general, these approaches use canonical solutions to train the models and often find a student’s code snapshot in a model code progression to determine what a suggested next step might be. While superficially similar to our approach, usage of our system differs in two fundamental ways. The first is that ITSs generally suggest a next step, whereas our feedback is potentially retrospective. Second, we are interested in both step-wise and holistic feedback. The idea of commentary on the overall arc of developing a program is not a prominent concept in ITSs.

Programming process data is a largely untapped resource in understanding how students write code. This is recognized by researchers who have proposed questions that could be answered given a suitable feature set and analysis techniques. Shrestha et al. [80] asked what exactly good programming process is. Is it student dependent? Is it dependent on the type of problem? If so, how? Edwards et al. [20] asked what, exactly, flailing and learning looks like. Some work on this has been done [81] using time spent on exercises and submission frequency as features. We suggest that richer features would lead to deeper understanding of student cognition and behavior.

3 Possibilities of Large Language Models for Programming Process Data

Large language models provide many exciting opportunities for analyzing and utilizing programming process data. One potential application of LLMs is to summarize the process into natural language. Before LLMs, most work in programming process data used derived metrics like time-on-task, number of events, etc. to analyze the process [33, 46], but here the utility of the data is only as good as the metrics derived from it. With LLMs, more abstract aspects of the process could potentially be derived. Another stream of work has looked into visualizing the programming process based on programming process data [19, 29, 80, 88], but here one issue is that the visualizations might be hard to interpret, or if the whole process is visualized keystroke-by-keystroke (see e.g. [19, 29, 88]), it is very time consuming to analyze.

Another potential application of LLMs for programming process data is to give feedback to students. When the whole process is available, feedback could be given with different granularities. For example, it could be possible to give high-level feedback on the whole process (e.g., “*you should run code more frequently*”) or low-level feedback that considers the process (e.g., “*you recently modified your function ‘calculateSum’ but there is a small bug*”). Ideally, feedback using LLMs could be given *during* the process, which could support students while they are working on their programs.

Some recent work has looked into generating synthetic data using LLMs [26, 49, 58, 68, 69], which is one potential application of LLMs for programming process data. One issue with programming process data, especially very fine-grained data such as keystroke data, is that individuals can sometimes be identified using “keystroke dynamics”, i.e., based on their unique typing rhythms [3, 57, 82]. This is an issue for sharing data openly [20, 52], since typically data that contains student personal information or

¹<https://pslcdatashop.web.cmu.edu/Files?datasetId=3458>

that can be traced back to the student cannot be shared. In Europe, the General Data Protection Regulation (GDPR) specifically lists biometric data (which keystroke data falls under [82]) under “special categories of personal data” which have more strict requirements under the GDPR. While some work has looked into the de-identification of keystroke data [20, 52], de-identification can reduce the utility of the data [52]. Thus, generating synthetic programming process data from authentic programming process data utilizing LLMs could help advance research, as long as the generated data is equal to real data quality-wise.

One challenge in utilizing process data with LLMs is the cost of using LLMs, which for many commercial LLMs is based on the number of tokens used². As process data is typically large [20, 33], the cost could be great. However, as the context windows of LLMs continue to increase³, and the capabilities of open-source models increase, the costs of using LLMs might reduce.

4 Case Study: Analyzing Programming Process Data with LLMs

Here, we outline our case study of analyzing programming process data using LLMs. For the case study, we explored how well three different LLMs 1) summarize the programming process based on the programming process data and 2) generate programming process feedback for students based on the programming process data.

4.1 Process Data and Preprocessing

For the process data, we used a publicly available Python dataset published by Edwards et al. [20]. The dataset has 44 participants and 8 assignments from an introductory programming course organized at Utah State University in the United States. The dataset was collected using the PyPhanon plugin [21] for PyCharm⁴, and it consists of 1 million unique events. For the present work, we focused on three of the assignments. Outlines of the assignments are given in Table 1 and are taken from assignment descriptions included with the dataset.

Using the logs from five students working on three assignments, we construed the programming process from the logs into a sequence of state snapshots that outlined how the assignment code evolved over time from the beginning to the end. For the present analysis, we created the sequence of state snapshots by including the first and last code states and each code snapshot immediately preceding breaks of five minutes or more. Table 2 provides descriptive statistics on the number of state snapshots for each of the students for the included assignments.

4.2 Large Language Models and Prompt Engineering

For the analysis, we used Anthropic’s Claude 3 Opus⁵, OpenAI’s GPT-4 Turbo⁶, and Meta’s LLaMa2 70B Chat model⁷. The first two models are proprietary, while the third model is an open-source model. The models were selected to represent a snapshot of the state-of-the-art LLMs at the time of the analysis in March and April of 2024. We acknowledge that OpenAI has the GPT-4 model that is known to perform better than the GPT-4 Turbo. However, the context window size for the GPT-4 version that was available for us was limited to 8k tokens, and our analyses highlighted that this was not sufficient. Thus, we opted for the GPT-4 Turbo which has a context window size of 128k tokens. This is also a limitation in some of our analyses where we leverage the LLaMa2 70B Chat model.

Our prompt engineering had the objective of having LLMs describe the process a student took to write their computer program and then to provide feedback to the student as if the AI were an expert teaching assistant (TA). Following prompting best practices [22], we iteratively explored and refined a range of prompts, including providing context and domain information, personalizing the responses, providing information about the format of the input data, and providing guidelines and constraints on the expected outcomes. One of the authors was in charge of the prompt engineering process, sharing intermediate results and observations with the other authors. During the prompt engineering process, the research team met weekly, commenting on the outputs and the prompts, and discussing further prompt improvement possibilities and briefly evaluating their impact on the outcomes.

During the prompt engineering process, we observed that the models were relatively poor at incorporating timestamp information. As an example, if a student had a break that consisted of multiple days, an LLM might suggest that the student was thinking hard about the problem for days. Thus, we omitted the timestamp data. We further explored specifying and not specifying the number of items in the feedback, and in the end resorted to not restricting the number of items for the present evaluation.

The final prompt for providing feedback on the programming process is outlined in Figure 1. When comparing the prompt for providing feedback on the programming process to the prompt for summarizing the programming process, the main differences are in the persona “[...] and an expert in summarizing how students construct their programs [...]” and the explicit task instructions “As an extremely good introductory programming teaching assistant, summarize how the student constructed their program based on the following programming process data. Respond as if you were talking to the student. Only summarize the programming process data. Do not provide any suggestions or feedback.”.

²See, for example, OpenAI’s pricing: <https://openai.com/api/pricing/> (Retrieved 4 June 2024).

³For example, Google’s Gemini has a context window of up to ten million tokens: <https://blog.google/technology/ai/google-gemini-next-generation-model-february-2024/#context-window> (Retrieved 4 June 2024).

⁴<https://www.jetbrains.com/pycharm/>

⁵<https://www.anthropic.com/news/claude-3-family>,
claude-3-opus-20240229

version

⁶<https://platform.openai.com/docs/models/gpt-4-and-gpt-4-turbo>,
gpt-4-0125-preview

version

⁷<https://llama.meta.com/llama2>,
meta-llama/Llama-2-70b-chat-hf

HuggingFace

version

Table 1: Summary descriptions of the three assignments included in the case study.

Assignment name	Brief description
Fluky numbers	Calculate fluky numbers. “You will write a program that will find Fluky Numbers. A Fluky Number is a number that is equal to a random number, where the random number is the first random number generated after seeding a random number generator with the sum of all factors of a number, not including itself.”
Zookeeper	Zookeeper and elephants in pens. “Your job is to write a program to simulate [a zookeeper checking elephant pens] 100,000 times to determine [the probabilities of the zookeeper finding an elephant in each pen].”
Number pyramids	Print a number pyramid. “The user will enter the number of rows to have in a pyramid of numbers. Based on the user’s input the program will display a pyramid where each row contains the number x , which is printed x times, where x is the row number.”

Table 2: Summary of number of states in the programming process data included in the data for LLMs for analysis.

Student	Average states	Fluky numbers states	Zookeeper states	Number pyramid states
S1	9.3	16	8	5
S2	8	11	6	7
S3	7.7	10	6	7
S4	7.7	9	3	11
S5	5.3	6	4	6
Avg.	7.6	10.4	5.4	7.2

4.3 Data Generation

To generate the data, we used Anthropic’s and OpenAI’s APIs for Claude 3 Opus and GPT-4 Turbo respectively, while LLaMa2 70B Chat model was run on the HuggingFace platform⁸. The data generation led to a total of 45 programming process descriptions (3 assignments \times 5 students \times 3 models) and 45 feedbacks based on the programming process, leading to 90 entries. Table 3 provides descriptive statistics of the data generation process, outlining the average response time and the average response length for each of the models.

4.4 Evaluation

The evaluation was divided into three parts. First, one of the researchers conducted a surface-level summary evaluation of the outputs to assess whether the models followed the task. Next, three researchers analyzed a subset of the data to develop a shared understanding, which was discussed collectively. Finally, informed by this shared understanding, one researcher conducted a thematic analysis of the improvement suggestions in the feedback to identify recurring themes in the outputs.

The thematic analysis was exploratory. The researcher identified recurring themes in the outputs, merging similar themes where appropriate. This approach was considered acceptable due to the open-ended nature of the problem and the length of the responses produced by the LLMs. The themes identified – constructive and positive feedback, back-and-forth coding, implementation not aligned with the problem statement, referring to the problem statement, incremental implementation and testing, planning ahead,

reducing unnecessary outputs, removing unnecessary comments, dividing code into functions, variable and function naming, and hallucination – are discussed with examples in Section 4.5.2. Section 4.5.3 synthesizes the researchers’ discussions while forming the shared understanding of the data.

Detailed mappings of themes to feedback were not maintained due to the exploratory nature of the analysis. As a result, statistics on the presence of the themes are not available for the current study.

4.5 Results

4.5.1 Surface-Level Analysis. On a surface level, the LLMs followed the instructions and provided expected outputs adequately. For the 15 inputs (5 students and 3 assignments), GPT-4 Turbo generated 15 acceptable (see below for our definition of acceptable) programming process summaries and 14 feedbacks, Claude 3 Opus generated 13 acceptable summaries and 11 feedbacks, and LLaMa2 70B Chat generated 10 acceptable summaries and 7 feedbacks. While GPT-4 Turbo and Claude 3 Opus provided a response for all student-exercise-prompt combinations, LLaMa 70B Chat failed to provide a response in two instances (one summary and one feedback) due to the prompt being too long.

The main reason why a summary was not deemed acceptable was that the response focused on describing one of the code states but did not summarize the process. Similarly, the main reason when a programming process feedback was not deemed acceptable was that the feedback did not explicitly include information about the process, but provided more generic improvement suggestions such as suggesting more efficient ways to implement the algorithms, adding type hints, and so on.

⁸<https://huggingface.co/>

Table 3: The average response time and the average response length for each of the models.

Model	Average response time (seconds)	Average response length (characters)
Claude 3 Opus	28.6	1814
GPT-4 Turbo	30.2	3086
LLaMa2 70B Chat	14.2	1943

You are an introductory programming teaching assistant who is an expert in analyzing programming process data and an expert in providing suggestions on improving the programming process based on the programming process data (1). You are studying how a student solved a programming problem by analyzing the programming process data (2).

Here's the handout for the programming problem:

<handout> (3)

The programming process data is described as a sequence of steps, where each step is numbered and each step has the assignment header and the code. The format of the data is as follows: (4)

####

Step: step identifier

####

(the code state that the student had at this step)

As an extremely good introductory programming teaching assistant, provide suggestions on how to improve the programming process based on the following programming process data. Respond as if you were talking to the student. Only provide improvement suggestions based on the data. (5)

<process data> (6)

Figure 1: Prompt format for providing feedback on the programming process. The prompt included (1) a description of the context and the persona, (2) the broader task, (3) the handout, (4) the format of the input data, (5) the explicit task instructions, and (6) the process data.

4.5.2 Analysis of Programming Process Improvement Suggestions.

Overall, while the models did often provide similar suggestions, there were also differences in the focus and style of the feedback. The GPT-4 Turbo model often used a logical structure to the feedback consisting of generic comments, specific comments, and a conclusion. The feedbacks from Claude 3 Opus and LLaMa 70B Chat model on the other hand were typically more list-like, providing suggestions directly. When considering the programming process feedback as a whole, the models typically followed the implicit instruction of being constructive and providing feedback in a positive fashion as the feedback was directed to students. The feedback items often included something that the student was praised on, even if just for adding comments or for exploring, as shown in the following quotes (mainly from GPT-4 Turbo).

"It's great to see you organizing your code with comments, making it easier to understand the purpose of each section. Keep maintaining this practice."

"It's great that you started by experimenting with the random number generation in steps 1 and 3 to understand how it works. This exploratory coding is a good way to familiarize yourself with new concepts."

The models also often captured back and forth movement where students first added code and then removed it, as well as scenarios where students started working on implementations that were not in line with the problem statement. Such feedbacks included comments that would suggest the student to step away from coding to thinking about the problem, as shown by the following quotes.

"The initial steps indicate a bit of confusion regarding the problem statement. [...] Taking time to fully understand the problem before starting to code can save a lot of time and confusion."

"You seem to jump into implementing parts of the solution without fully forming the algorithm in your mind or on paper first. It's important to sketch out a basic algorithm or pseudocode before you begin coding, especially for a problem with several specific steps like this one."

The feedbacks also at times referenced to the assignment handout that was included in the prompt, as shown in the next quote. The handout for the "Fluky numbers" assignment had explicit hints on how to efficiently solve the problem, and the student had followed the hint in their implementation.

"You had the right idea in step 008 to keep track of the number of Fluky Numbers found and to terminate the search after finding seven. This directly aligns with the hint provided and is crucial for optimizing performance."

The feedbacks also included more common aspects that could be also non-assignment -specific. These suggestions included incremental implementation and testing of the solution, planning ahead, reducing unnecessary outputs and comments, removing commented out code, adding comments to code to improve readability, dividing code into functions, and other aspects such as variable and function naming. The following example starts by orienting the student with specific references to code snapshots and then uses that to make the general suggestion of not leaving commented-out code in the final submission. It then explains why the student should follow the suggestion (for code readability).

"As you iterated from Step 002 to Step 005, it looks like you attempted to address the issue of centering the pyramid but struggled to find the right solution on your first

tries. This is completely normal and part of the learning process. However, it's important to not leave commented-out code or unused variables in your final submission, as it can make the code harder to read. Find the solution that works best and then refactor your code to remove any unnecessary parts."

In one submission, the student used an obscenity in a comment. The feedback included this:

*"In Step 001, you've used an unfortunate typo, 'fu**tion', which seems to be an attempt at humor but can come across as unprofessional, especially in a learning or professional setting. In the future, I recommend keeping comments and code professional and clear. This will not only help maintain a positive learning environment but also prepare you for professional coding practices."*

The LLMs in the context of our study were not immune to hallucinations. For example, the following feedback recommends one function instead of another for performance reasons when, in fact, `random.randint(0, k)` is an alias for `random.randrange(0, k+1)`:

"The code uses the 'random.randint(0, k)' function to generate a random number between 0 and 'k'. However, this function can be slow for large values of 'k'. Instead, use the 'random.randrange(0, k+1)' function, which is faster and more efficient."

4.5.3 Researcher Insights From Programming Process Summarization and Feedback. The AI systems were explicitly instructed not to include feedback on the code itself, but rather to focus on the process of writing the code. (One round of prompting includes the statement, "Only provide feedback on how to improve the programming process. Do not provide any summary description of the programming process data and do not provide explicit suggestions on how to improve the code." Other prompts include similar statements.) However, the LLMs were found to have great difficulty in recognizing the difference between the process a student takes to write code and the code itself. For example, the AI gave the following process-related feedback to one student:

"Incremental Development and Testing: I noticed that you began writing large blocks of code before testing smaller units. It's important to approach complex problems by breaking them down into smaller, more manageable pieces. Once you write a small part, such as a function to calculate the sum of factors, test it thoroughly before moving on. This will help you catch errors early and understand each part of your code better. Consider implementing your code piece by piece, testing each function independently before integrating it into your main logic."

The same feedback set also includes the following, which relates to the code itself:

"Efficiency in Your Loops and Functions: Your approach to finding the Fluky Numbers involves quite a bit of redundancy and potentially unnecessary calculations. For example, in calculating the sum of factors, think

about how you might optimize this process to avoid redundant calculations or checks. In addition, I see you've modified the 'factorSum' function to return both a sum and a count of factors, but the count doesn't appear to be utilized in later logic. Always review your functions to ensure they are doing exactly what is needed for the task at hand, no more, no less. This will help keep your code efficient and focused."

The fact that the feedback includes the word "process" in the feedback indicates that it may be confusing the process of code execution with the process of code writing. Every single summary and feedback output from the LLMs included something relating to the code and not process. In retrospect, this should not have been surprising to us. Even we humans tend to find it challenging to distinguish between the two. One possible reason for this is that software engineers and programming educators are attuned to thinking about code but think less about the evolution of code into its final form. Thus, when we look at a code snapshot, we tend to analyze it as a standalone product, rather than in its temporal context. Furthermore, considering it through the lens of cognitive load theory, analyzing a code snapshot relative to a prior snapshot requires roughly double the working memory, which is already stretched when reading code. From a philosophical perspective, it makes sense that LLMs suffer from this as well, as they are trained on data generally generated by humans, and there is very little data that deals with the programming process.

Furthermore, the majority of generated summaries and feedback that are related to the programming process are generic. The above example that suggests that the student use incremental development and testing is typical: incremental development is so commonplace that the only thing keeping it from becoming cliché is that it is so rarely actually practiced. Whereas our vocabulary of things that can be improved in code is vast (e.g., variable names, comments, duplicated code, algorithm complexity, code organization into functions, etc.), educators have far fewer candidate improvements to the programming process, which are generally limited to incremental development, early decomposition, and planning ahead.

A challenge in evaluating quality of summaries and feedback is their subjectivity. Feedback in particular, is a challenge because, in some cases, educators disagree on best practices. For example, many educators encourage up-front design whereas others encourage a more agile, refactoring approach. In order to make evaluation of AI output objective, the AI would need to be primed with opinions on controversial topics. Related to this, educators would need to decide which topics are the most important. In our case study, the AI responses were often excessively long, so the educator would need to limit the length of output and indicate to the AI which topics to prioritize.

So what criteria should an AI be judged on? We propose five scores. The first would be a hallucination score. This would be straightforward to measure with standard multi-rater coding practices. The second score would measure how well the AI distinguishes between process and analysis of static code. As discussed, this distinction is challenging for both LLMs and humans. However, we found that determining how well the AI makes the distinction is actually not too difficult. The human evaluator needn't consult the

code snapshots to determine if the LLM output addresses process or code. The third score would be specificity against genericness relative to the given submission. An AI that speaks in tropes is not useful. And finally, scores on correctness and utility would be needed. These are straightforward measures to reason about but are actually the most difficult to evaluate. It requires the evaluator to consider code snapshots relative to other snapshots which, as discussed above, is challenging.

One last note: we found that the AI systems distinguish extremely well between generating summaries and generating feedback. That is, when we ask the LLM to give a summary, the output never focuses on feedback, and vice versa.

5 Discussion

Overall, we here make an argument that exploring the ability of LLMs to provide feedback on programming process is worthwhile. The argument is in two parts. The first is the argument that feedback itself on programming process is useful. While giving students feedback on code is common, giving them feedback on how they went about writing the code is rare. This is because process data isn't easy to collect, but also because little research has been done into best practices in code evolution in the context of novice programmers. We claim that guiding students on this front will result in less frustration, better learning, reduced attrition, and increased diversity in the computer programming community. If we accept that providing guidance to students on process is useful, then we can address the other part of the argument, that AI in general, and LLMs specifically, are well-suited to providing this feedback to students. The primary reason for this is that analyzing how a code evolved is far more difficult than analyzing a single code snapshot. Well-prompted LLMs have vast memory resources and can evaluate far more code submissions than a human, and can do it in a short amount of time.

One way to make process feedback even more effective could be to combine it with a playback tool such as *KeystrokeExplorer* [20]. We note that in our case study, we considered only snapshots of code immediately preceding a break of at least five minutes. During our analysis, we found that this granularity often left the reader disoriented. Actually watching the keystroke-level evolution of the code between snapshots could be helpful in keeping the student oriented as they receive the AI-generated feedback, and also potentially create self-reflection opportunities, in a similar way as one could, e.g., review their chess games.

We see that already finding out the appropriate granularity (and data format) could significantly improve the preliminary results outlined in this article, not to mention the emergence of newer and more powerful LLMs. We also see that using LLMs for creating feedback on the programming process data and providing that feedback to students during the programming process could help in pinning down the points in time when feedback should be given, which has been an open question in prior programming process feedback research [35]. To begin the discussion at the conference, and beyond, we outline the following five research directions for using LLMs on programming process data:

- Cost-efficient representations of programming process data for LLMs, including evaluating data representation types

such as diffs and exploring the utility of different data granularities.

- Providing insights from programming process data to students, exploring when and how to present the insights for maximal effectiveness, including incorporating the LLM-generated insights into existing systems such as IDEs and programming process visualizers.
- Analyzing the subjectivity of LLM-generated insights and their utility to learners and instructors, and identifying contextual factors that contribute to the subjectivity.
- Exploring the opportunities of leveraging programming process data with LLMs for supporting students in learning about the process of programming, which has been classically highlighted as one of the challenges in learning complex skills such as programming [7, 16, 84].
- Training and fine-tuning LLMs with processing programming process data for more efficient and higher quality programming process insights.

Acknowledgments

This research was supported by the Research Council of Finland (Academy Research Fellow grant number 356114).

References

- [1] Alireza Ahadi, Raymond Lister, Heikki Haapala, and Arto Vihavainen. 2015. Exploring machine learning methods to automatically identify students in need of assistance. In *ICER '15: Proceedings of the Eleventh Annual International Conference on International Computing Education Research*, Judy Sheard and Quintin Cutts (Eds.). Association for Computing Machinery, 121–130.
- [2] Alireza Ahadi, Raymond Lister, and Arto Vihavainen. 2016. On the number of attempts students made on some online programming exercises during semester and their subsequent performance on final exam questions. In *ITiCSE '16: Proceedings of the 2016 ACM Conference on Innovation and Technology in Computer Science Education*. Association for Computing Machinery, 218–223.
- [3] Md Liakat Ali, John V Monaco, Charles C Tappert, and Meikang Qiu. 2017. Keystroke biometric systems for user authentication. *Journal of Signal Processing Systems* 86 (2017), 175–190.
- [4] Tiffany Barnes and John Stamper. 2008. Toward automatic hint generation for logic proof tutoring using historical student data. In *International Conference on Intelligent Tutoring Systems*, Beverley P. Woolf, Esma Aïmeur, Roger Nkambou, and Susanne Lajoie (Eds.). Springer, 373–382.
- [5] Brett A Becker. 2016. A new metric to quantify repeated compiler errors for novice programmers. In *ITiCSE '16: Proceedings of the 2016 ACM Conference on Innovation and Technology in Computer Science Education*. Association for Computing Machinery, 296–301.
- [6] Brett A Becker, Paul Denny, Raymond Pettit, Durell Bouchard, Dennis J Bouvier, Brian Harrington, Amir Kamil, Amey Karkare, Chris McDonald, Peter-Michael Osera, et al. 2019. Compiler error messages considered unhelpful: The landscape of text-based programming error message research. *Proceedings of the working group reports on innovation and technology in computer science education* (2019), 177–210.
- [7] Jens Bennedsen and Michael E Caspersen. 2005. Revealing the programming process. In *Proceedings of the 36th SIGCSE technical symposium on Computer science education*. 186–190.
- [8] Paulo Blikstein, Marcelo Worsley, Chris Piech, Mehran Sahami, Steven Cooper, and Daphne Koller. 2014. Programming pluralism: Using learning analytics to detect patterns in the learning of computer programming. *Journal of the Learning Sciences* 23, 4 (2014), 561–599.
- [9] Neil CC Brown, Amjad Altadmri, Sue Sentance, and Michael Kölling. 2018. Blackbox, five years on: An evaluation of a large-scale programming data collection project. In *ICER '18: Proceedings of the 2018 ACM Conference on International Computing Education Research*. Association for Computing Machinery, 196–204.
- [10] Neil CC Brown and Michael Kölling. 2020. Blackbox Mini-Getting Started With Blackbox Data Analysis. In *SIGCSE '20: Proceedings of the 51st ACM Technical Symposium on Computer Science Education*. Association for Computing Machinery, 1387–1387.
- [11] Neil CC Brown, Pierre Weill-Tessier, Juho Leinonen, Paul Denny, and Michael Kölling. 2024. Howzat? Appealing to Expert Judgement for Evaluating Human

- and AI Next-Step Hints for Novice Programmers. *arXiv preprint arXiv:2411.18151* (2024).
- [12] Neil Christopher Charles Brown, Michael Kölling, Davin McCall, and Ian Utting. 2014. Blackbox: A large scale repository of novice programmers' activity. In *SIGCSE '14: Proceedings of the 45th ACM Technical Symposium on Computer Science Education*. Association for Computing Machinery, 223–228.
 - [13] Adam Scott Carter and Christopher David Hundhausen. 2017. Using programming process data to detect differences in students' patterns of programming. In *Proceedings of the 2017 ACM SIGCSE Technical Symposium on Computer Science Education*. Association for Computing Machinery, 105–110.
 - [14] Adam S Carter, Christopher D Hundhausen, and Olusola Adesope. 2015. The normalized programming state model: Predicting student performance in computing courses based on programming behavior. In *ICER '15: Proceedings of the Eleventh Annual International Conference on International Computing Education Research*. Association for Computing Machinery, 141–150.
 - [15] Karo Castro-Wunsch, Alireza Ahadi, and Andrew Petersen. 2017. Evaluating neural networks as a method for identifying students in need of assistance. In *SIGCSE '17: Proceedings of the 2017 ACM Technical Symposium on Computer Science Education*. Association for Computing Machinery, 111–116.
 - [16] Allan Collins, John Seely Brown, Ann Holum, et al. 1991. Cognitive apprenticeship: Making thinking visible. *American educator* 15, 3 (1991), 6–11.
 - [17] Paul Denny, Brett A Becker, Juho Leinonen, and James Prather. 2023. Chat Overflow: Artificially Intelligent Models for Computing Education—renAIssance or apocAlypse?. In *Proceedings of the 2023 Conference on Innovation and Technology in Computer Science Education V. 1*. 3–4.
 - [18] Paul Denny, James Prather, Brett A Becker, James Finnie-Ansley, Arto Hellas, Juho Leinonen, Andrew Luxton-Reilly, Brent N Reeves, Eddie Antonio Santos, and Sami Sarsa. 2024. Computing education in the era of generative AI. *Commun. ACM* 67, 2 (2024), 56–67.
 - [19] Joseph Dittton, Hillary Swanson, and John Edwards. 2021. External imagery in computer programming. In *Proceedings of the 52nd ACM Technical Symposium on Computer Science Education*. 1226–1231.
 - [20] John Edwards, Kaden Hart, and Raj Shrestha. 2023. Review of CSED Data and Introduction of Two Public CS1 Keystroke Datasets. *Journal of Educational Data Mining* 15, 1 (2023), 1–31.
 - [21] John Edwards, Kaden Hart, and Christopher Warren. 2022. A practical model of student engagement while programming. In *Proceedings of the 2022 ACM SIGCSE technical symposium on computer science education*.
 - [22] Sabit Ekin. 2023. Prompt engineering for ChatGPT: a quick guide to techniques, tips, and best practices. *Authorea Preprints* (2023).
 - [23] Anthony Estey and Yvonne Coady. 2016. Can interaction patterns with supplemental study tools predict outcomes in CS1?. In *ITICSE '16: Proceedings of the 2016 ACM Conference on Innovation and Technology in Computer Science Education*. Association for Computing Machinery, 236–241.
 - [24] Anthony Estey, Hieke Keuning, and Yvonne Coady. 2017. Automatically classifying students in need of support by detecting changes in programming behaviour. In *SIGCSE '17: Proceedings of the 2017 ACM Technical Symposium on Computer Science Education*. Association for Computing Machinery, 189–194.
 - [25] Elena L Glassman, Jeremy Scott, Rishabh Singh, Philip J Guo, and Robert C Miller. 2015. OverCode: Visualizing variation in student solutions to programming problems at scale. *ACM Transactions on Computer-Human Interaction (TOCHI)* 22, 2 (2015), 1–35.
 - [26] Perttu Hämäläinen, Mikke Tavast, and Anton Kunnari. 2023. Evaluating Large Language Models in Generating Synthetic HCI Research Data: a Case Study. In *Proceedings of the Conference on Human Factors in Computing Systems (CHI '23)*.
 - [27] Kaden Hart, Christopher Warren, and John Edwards. 2023. Accurate Estimation of Time-on-Task While Programming. In *ACM Technical Symposium on Computing Science Education (SIGCSE)*.
 - [28] John Hattie and Helen Timperley. 2007. The power of feedback. *Review of educational research* 77, 1 (2007), 81–112.
 - [29] Kenny Heinonen, Kasper Hirvikoski, Matti Luukkainen, and Arto Vihavainen. 2014. Using codebrowser to seek differences between novice programmers. In *Proceedings of the 45th ACM technical symposium on Computer science education*. 229–234.
 - [30] Arto Hellas, Juho Leinonen, and Petri Ihanntola. 2017. Plagiarism in take-home exams: help-seeking, collaboration, and systematic cheating. In *ITICSE '17: Proceedings of the 2017 ACM conference on innovation and technology in computer science education*. Association for Computing Machinery, 238–243.
 - [31] Arto Hellas, Juho Leinonen, Sami Sarsa, Charles Koutchme, Lilja Kujanpää, and Juha Sorva. 2023. Exploring the responses of large language models to beginner programmers' help requests. In *Proceedings of the 2023 ACM Conference on International Computing Education Research-Volume 1*. 93–105.
 - [32] Jack Hollingsworth. 1960. Automatic graders for programming classes. *Commun. ACM* 3, 10 (1960), 528–529.
 - [33] Petri Ihanntola, Arto Vihavainen, Alireza Ahadi, Matthew Butler, Jürgen Börstler, Stephen H. Edwards, Essi Isohanni, Ari Korhonen, Andrew Petersen, Kelly Rivers, Miguel Ángel Rubio, Judy Sheard, Bronius Skupas, Jaime Spacco, Claudia Szabo, and Daniel Toll. 2015. Educational Data Mining and Learning Analytics in Programming: Literature Review and Case Studies. In *Proceedings of the 2015 ITICSE on Working Group Reports (Vilnius, Lithuania) (ITICSE-WGR '15)*. Association for Computing Machinery, New York, NY, USA, 41–63. <https://doi.org/10.1145/2858796.2858798>
 - [34] Matthew C Judad. 2006. Methods and tools for exploring novice compilation behaviour. In *ICER '06: Proceedings of the second International Workshop on Computing Education Research*. Association for Computing Machinery, 73–84.
 - [35] Johan Jeuring, Hieke Keuning, Sami Sarsa, Dennis Bouvier, Cruz Izu, Natalie Kiesler, Teemu Lehtinen, Dominic Lohr, Andrew Peterson, and Sami Sarsa. 2022. Towards giving timely formative feedback and hints to novice programmers. In *Proceedings of the 2022 Working Group Reports on Innovation and Technology in Computer Science Education*. 95–115.
 - [36] Ayaan M Kazerouni, Stephen H Edwards, T Simin Hall, and Clifford A Shaffer. 2017. DevEventTracker: Tracking development events to assess incremental development and procrastination. In *ITICSE '17: Proceedings of the 2017 ACM Conference on Innovation and Technology in Computer Science Education*. Association for Computing Machinery, 104–109.
 - [37] Ayaan M Kazerouni, Stephen H Edwards, and Clifford A Shaffer. 2017. Quantifying incremental development practices and their relationship to procrastination. In *Proceedings of the 2017 ACM Conference on International Computing Education Research*, Josh Tenenber and Lauri Malmi (Eds.). Association for Computing Machinery, 191–199.
 - [38] Hieke Keuning, Johan Jeuring, and Bastiaan Heeren. 2018. A Systematic Literature Review of Automated Feedback Generation for Programming Exercises. *ACM Trans. Comput. Educ.* 19, 1, Article 3 (2018), 43 pages.
 - [39] Teemu Koivisto and Arto Hellas. 2022. Evaluating CodeClusters for Effectively Providing Feedback on Code Submissions. In *2022 IEEE Frontiers in Education Conference (FIE)*. IEEE, 1–9.
 - [40] Charles Koutchme, Nicola Dainese, Sami Sarsa, Arto Hellas, Juho Leinonen, and Paul Denny. 2024. Open Source Language Models Can Provide Feedback: Evaluating LLMs' Ability to Help Students Using GPT-4-As-A-Judge. *arXiv preprint arXiv:2405.05253* (2024).
 - [41] Charles Koutchme, Sami Sarsa, Arto Hellas, Lassi Haaranen, and Juho Leinonen. 2022. Methodological Considerations for Predicting At-risk Students. In *Australasian Computing Education Conference*. Association for Computing Machinery, 105–113.
 - [42] Charles Koutchme, Sami Sarsa, Juho Leinonen, Arto Hellas, and Paul Denny. 2023. Automated program repair using generative models for code infilling. In *International Conference on Artificial Intelligence in Education*. Springer, 798–803.
 - [43] Sowndarya Krishnamoorthy, Luis Rueda, Sherif Saad, and Haytham Elmiligi. 2018. Identification of user behavioral biometrics for authentication using keystroke dynamics and machine learning. In *ICBEA '18: Proceedings of the 2018 2nd International Conference on Biometric Engineering and Applications*. Association for Computing Machinery, 50–57.
 - [44] Teemu Lehtinen, Lassi Haaranen, and Juho Leinonen. 2023. Automated questionnaires about students' JavaScript programs: towards gauging novice programming processes. In *Proceedings of the 25th Australasian Computing Education Conference*. 49–58.
 - [45] Antti Leinonen, Henrik Nygren, Nea Pirttinen, Arto Hellas, and Juho Leinonen. 2019. Exploring the applicability of simple syntax writing practice for learning programming. In *SIGCSE '19: Proceedings of the 50th ACM Technical Symposium on Computer Science Education*. Association for Computing Machinery, 84–90.
 - [46] Juho Leinonen. 2019. *Keystroke Data in Programming Courses*. Ph.D. Dissertation. University of Helsinki.
 - [47] Juho Leinonen, Francisco Enrique Vicente Castro, and Arto Hellas. 2022. Time-on-Task Metrics for Predicting Performance. In *SIGCSE '22: Proceedings of the 53rd ACM Technical Symposium on Computer Science Education*. Association for Computing Machinery, 871–877.
 - [48] Juho Leinonen, Francisco Enrique Vicente Castro, Arto Hellas, et al. 2021. Fine-Grained Versus Coarse-Grained Data for Estimating Time-on-Task in Learning Programming. In *Proceedings of The 14th International Conference on Educational Data Mining (EDM 2021)*, Sharon Hsiao, Shaghayegh Sahebi, Francois Bouchet, and Jill-J'enn Vie (Eds.). The International Educational Data Mining Society, 648–653.
 - [49] Juho Leinonen, Paul Denny, Olli Kiljunen, Stephen MacNeil, Sami Sarsa, and Arto Hellas. 2024. LLM-itation is the Sincerest Form of Data: Generating Synthetic Buggy Code Submissions for Computing Education. *arXiv preprint arXiv:2411.10455* (2024).
 - [50] Juho Leinonen, Paul Denny, Stephen MacNeil, Sami Sarsa, Seth Bernstein, Joanne Kim, Andrew Tran, and Arto Hellas. 2023. Comparing code explanations created by students and large language models. In *Proceedings of the 2023 Conference on Innovation and Technology in Computer Science Education V. 1*. 124–130.
 - [51] Juho Leinonen, Arto Hellas, Sami Sarsa, Brent Reeves, Paul Denny, James Prather, and Brett A Becker. 2023. Using large language models to enhance programming error messages. In *Proceedings of the 54th ACM Technical Symposium on Computer Science Education V. 1*. 563–569.
 - [52] Juho Leinonen, Petri Ihanntola, and Arto Hellas. 2017. Preventing Keystroke Based Identification in Open Data Sets. In *Proceedings of the Fourth (2017) ACM*

- Conference on Learning @ Scale* (Cambridge, Massachusetts, USA) (*L@S '17*). Association for Computing Machinery, New York, NY, USA, 101–109. <https://doi.org/10.1145/3051457.3051458>
- [53] Juho Leinonen, Leo Leppänen, Petri Ihanntola, and Arto Hellas. 2017. Comparison of time metrics in programming. In *ICER '17: Proceedings of the 2017 ACM Conference on International Computing Education Research*. Association for Computing Machinery, 200–208.
- [54] Juho Leinonen, Krista Longi, Arto Klami, Alireza Ahadi, and Arto Vihavainen. 2016. Typing patterns and authentication in practical programming exams. In *ITiCSE '16: Proceedings of the 2016 ACM Conference on Innovation and Technology in Computer Science Education*. Association for Computing Machinery, 160–165.
- [55] Juho Leinonen, Krista Longi, Arto Klami, and Arto Vihavainen. 2016. Automatic inference of programming performance and experience from typing patterns. In *SIGCSE '16: Proceedings of the 47th ACM Technical Symposium on Computing Science Education*. Association for Computing Machinery, 132–137.
- [56] Evanfiya Logacheva, Arto Hellas, James Prather, Sami Sarsa, and Juho Leinonen. 2024. Evaluating Contextually Personalized Programming Exercises Created with Generative AI. *arXiv preprint arXiv:2407.11994* (2024).
- [57] Krista Longi, Juho Leinonen, Henrik Nygren, Joni Salmi, Arto Klami, and Arto Vihavainen. 2015. Identification of programmers from typing patterns. In *Proceedings of the 15th Koli Calling Conference on Computing Education Research* (Koli, Finland) (*Koli Calling '15*). Association for Computing Machinery, New York, NY, USA, 60–67. <https://doi.org/10.1145/2828959.2828960>
- [58] Stephen MacNeil, Magdalena Rogalska, Juho Leinonen, Paul Denny, Arto Hellas, and Xandria Crosland. 2024. Synthetic Students: A Comparative Study of Bug Distribution Between Large Language Models and Computing Students. In *Proceedings of the 2024 on ACM Virtual Global Computing Education Conference V. 1* (Virtual Event, NC, USA) (*SIGCSE Virtual 2024*). Association for Computing Machinery, New York, NY, USA, 137–143. <https://doi.org/10.1145/3649165.3690100>
- [59] Stephen MacNeil, Andrew Tran, Arto Hellas, Joanne Kim, Sami Sarsa, Paul Denny, Seth Bernstein, and Juho Leinonen. 2023. Experiences from using code explanations generated by large language models in a web software development e-book. In *Proceedings of the 54th ACM Technical Symposium on Computer Science Education V. 1*. 931–937.
- [60] Stephen MacNeil, Andrew Tran, Dan Mogil, Seth Bernstein, Erin Ross, and Ziheng Huang. 2022. Generating diverse code explanations using the gpt-3 large language model. In *Proceedings of the 2022 ACM Conference on International Computing Education Research-Volume 2*. 37–39.
- [61] Jessica McBroom, Irena Koprińska, and Kalina Yacef. 2021. A survey of automated programming hint generation: The hints framework. *ACM Computing Surveys (CSUR)* 54, 8 (2021), 1–27.
- [62] Aythami Morales and Julian Fierrez. 2015. Keystroke biometrics for student authentication: A case study. In *ITiCSE '15: Proceedings of the 2015 ACM Conference on Innovation and Technology in Computer Science Education*. Association for Computing Machinery, 337–337.
- [63] Christian Murphy, Gail Kaiser, Kristin Loveland, and Sahar Hasan. 2009. Retina: helping students and instructors based on observed programming activities. In *Proceedings of the 40th ACM technical symposium on Computer Science Education*, Gary Lewandowski and Steven Wolfman (Eds.). Association for Computing Machinery, 178–182.
- [64] Andy Nguyen, Christopher Piech, Jonathan Huang, and Leonidas Guibas. 2014. Codewebs: scalable homework search for massive open online programming courses. In *Proceedings of the 23rd international conference on World wide web*. 491–502.
- [65] Claudia Ott, Anthony Robins, and Kerry Shephard. 2016. Translating principles of effective feedback for students into the CS1 context. *ACM Transactions on Computing Education (TOCE)* 16, 1 (2016), 1–27.
- [66] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems* 35 (2022), 27730–27744.
- [67] José Carlos Paiva, José Paulo Leal, and Álvaro Figueira. 2022. Automated Assessment in Computer Science Education: A State-of-the-Art Review. *ACM Trans. Comput. Educ.* 22, 3, Article 34 (2022), 40 pages.
- [68] Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*. 1–22.
- [69] Joon Sung Park, Lindsay Popowski, Carrie Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2022. Social simulacra: Creating populated prototypes for social computing systems. In *Proceedings of the 35th Annual ACM Symposium on User Interface Software and Technology*. 1–18.
- [70] Chris Piech, Mehran Sahami, Daphne Koller, Steve Cooper, and Paulo Blikstein. 2012. Modeling how students learn to program. In *SIGCSE '12: Proceedings of the 43rd ACM Technical Symposium on Computer Science Education*. Association for Computing Machinery, 153–160.
- [71] James Prather, Paul Denny, Juho Leinonen, Brett A Becker, Ibrahim Albluwi, Michelle Craig, Hieke Keuning, Natalie Kiesler, Tobias Kohn, Andrew Luxton-Reilly, et al. 2023. The robots are here: Navigating the generative ai revolution in computing education. In *Proceedings of the 2023 Working Group Reports on Innovation and Technology in Computer Science Education*. 108–159.
- [72] Thomas Price, Rui Zhi, and Tiffany Barnes. 2017. Evaluation of a Data-Driven Feedback Algorithm for Open-Ended Programming. In *Proceedings of The 10th International Conference on Educational Data Mining (EDM 2017)*, X. Hu, T. Barnes, A. Hershkovitz, and L. Paquette (Eds.). International Educational Data Mining Society, 192–197.
- [73] Thomas W Price, Neil CC Brown, Dragan Lipovac, Tiffany Barnes, and Michael Kölling. 2016. Evaluation of a frame-based programming editor. In *Proceedings of the 2016 ACM Conference on International computing education research*. Association for Computing Machinery, 33–42.
- [74] Thomas W Price, Yihuan Dong, and Dragan Lipovac. 2017. iSnap: towards intelligent tutoring in novice programming environments. In *SIGCSE '17: Proceedings of the 2017 ACM Technical Symposium on Computer Science Education*. Association for Computing Machinery, 483–488.
- [75] Zac Pullar-Strecker, Filipe Dwan Pereira, Paul Denny, Andrew Luxton-Reilly, and Juho Leinonen. 2023. G is for Generalisation: Predicting Student Success from Keystrokes. In *Proceedings of the 54th ACM Technical Symposium on Computer Science Education V. 1*. 1028–1034.
- [76] Kelly Rivers and Kenneth R Koedinger. 2017. Data-driven hint generation in vast solution spaces: a self-improving python programming tutor. *International Journal of Artificial Intelligence in Education* 27, 1 (2017), 37–64.
- [77] Gustavo Rodriguez-Rivera, Jeff Turkstra, Jordan Buckmaster, Killian LeClainche, Shawn Montgomery, William Reed, Ryan Sullivan, and Jarett Lee. 2022. Tracking Large Class Projects in Real-Time Using Fine-Grained Source Control. In *SIGCSE '22: Proceedings of the 53rd ACM Technical Symposium on Computer Science Education*. Association for Computing Machinery, 565–570.
- [78] Eddie Antonio Santos, Prajish Prasad, and Brett A Becker. 2023. Always provide context: The effects of code context on programming error message enhancement. In *Proceedings of the ACM Conference on Global Computing Education Vol 1*. 147–153.
- [79] Sami Sarsa, Paul Denny, Arto Hellas, and Juho Leinonen. 2022. Automatic generation of programming exercises and code explanations using large language models. In *Proceedings of the 2022 ACM Conference on International Computing Education Research-Volume 1*. 27–43.
- [80] Raj Shrestha, Juho Leinonen, Arto Hellas, Petri Ihanntola, and John Edwards. 2022. CodeProcess Charts: Visualizing the Process of Writing Code. In *Proceedings of the 24th Australasian Computing Education Conference (Virtual Event, Australia) (ACE '22)*. Association for Computing Machinery, New York, NY, USA, 46–55. <https://doi.org/10.1145/3511861.3511867>
- [81] Jaime Spacco, Paul Denny, Brad Richards, David Babcock, David Hovemeyer, James Moscola, and Robert Duvall. 2015. Analyzing student work patterns using programming exercise data. In *SIGCSE '15: Proceedings of the 46th ACM Technical Symposium on Computer Science Education*. Association for Computing Machinery, 18–23.
- [82] Pin Shen Teh, Andrew Beng Jin Teoh, Shigang Yue, et al. 2013. A survey of keystroke dynamics biometrics. *The Scientific World Journal* 2013 (2013).
- [83] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).
- [84] Arto Vihavainen, Matti Paksula, and Matti Luukkainen. 2011. Extreme apprenticeship method in teaching programming for beginners. In *Proceedings of the 42nd ACM technical symposium on Computer science education*. 93–98.
- [85] Arto Vihavainen, Thomas Vikberg, Matti Luukkainen, and Martin Pärtel. 2013. Scaffolding students' learning using test my code. In *ITiCSE '13: Proceedings of the 18th ACM Conference on Innovation and Technology in Computer Science Education*. Association for Computing Machinery, 117–122.
- [86] Christopher Watson, Frederick WB Li, and Jamie I Godwin. 2013. Predicting performance in an introductory programming course by logging and analyzing student programming behavior. In *2013 IEEE 13th International Conference on Advanced Learning Technologies (ICALT)*. IEEE, 319–323.
- [87] Benedikt Wisniewski, Klaus Zierer, and John Hattie. 2020. The power of feedback revisited: A meta-analysis of educational feedback research. *Frontiers in psychology* 10 (2020), 487662.
- [88] Benjamin Xie, Jared Ordona Lim, Paul KD Pham, Min Li, and Amy J Ko. 2023. Developing Novice Programmers' Self-Regulation Skills with Code Replays. In *Proceedings of the 2023 ACM Conference on International Computing Education Research-Volume 1*. 298–313.
- [89] Lisa Yan, Annie Hu, and Chris Piech. 2019. Pensieve: Feedback on coding process for novices. In *SIGCSE '19: Proceedings of the 50th ACM Technical Symposium on Computer Science Education*. Association for Computing Machinery, 253–259.
- [90] Hongyan Zhong, Jun Niu, and Junjie Li. 2024. InsProg: Supporting Teaching Through Visual Analysis of Students' Programming Processes. In *Proceedings of the 2024 International Conference on Advanced Visual Interfaces*. 1–5.